

09/890929
PCT/JP00/08133
17.11.00

日 本 国 特 許 庁

PATENT OFFICE
JAPANESE GOVERNMENT

EU

JP00/8133

REC'D 19 JAN 2001

WIPO PCT

別紙添付の書類に記載されている事項は下記の出願書類に記載されて
いる事項と同一であることを証明する。

This is to certify that the annexed is a true copy of the following application as filed
with this Office.

出 願 年 月 日
Date of Application:

1999年12月14日

出 願 番 号
Application Number:

平成11年特許願第354401号

出 願 人
Applicant(s):

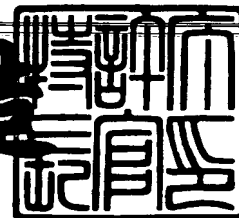
日立ソフトウェアエンジニアリング株式会社

PRIORITY
DOCUMENT
SUBMITTED OR TRANSMITTED IN
COMPLIANCE WITH RULE 17.1(a) OR (b)

2000年12月22日

特許庁長官
Commissioner,
Patent Office

及川耕造



出証番号 出証特2000-3105553

【書類名】 特許願

【整理番号】 11B016

【提出日】 平成11年12月14日

【あて先】 特許庁長官 殿

【国際特許分類】 G06F 17/00

【発明の名称】 樹状図表示方法及び樹状図表示システム

【請求項の数】 13

【発明者】

【住所又は居所】 神奈川県横浜市中区尾上町 6 丁目 8 1 番地 日立ソフト
ウェアエンジニアリング株式会社内

【氏名】 野崎 康行

【発明者】

【住所又は居所】 神奈川県横浜市中区尾上町 6 丁目 8 1 番地 日立ソフト
ウェアエンジニアリング株式会社内

【氏名】 渡辺 恒彦

【発明者】

【住所又は居所】 神奈川県横浜市中区尾上町 6 丁目 8 1 番地 日立ソフト
ウェアエンジニアリング株式会社内

【氏名】 中重 亮

【発明者】

【住所又は居所】 神奈川県横浜市中区尾上町 6 丁目 8 1 番地 日立ソフト
ウェアエンジニアリング株式会社内

【氏名】 田村 卓郎

【特許出願人】

【識別番号】 000233055

【氏名又は名称】 日立ソフトウェアエンジニアリング株式会社

【代理人】

【識別番号】 100091096

【弁理士】

【氏名又は名称】 平木 祐輔

【選任した代理人】

【識別番号】 100102576

【弁理士】

【氏名又は名称】 渡辺 敏章

【手数料の表示】

【予納台帳番号】 015244

【納付金額】 21,000円

【提出物件の目録】

【物件名】 明細書 1

【物件名】 図面 1

【物件名】 要約書 1

【包括委任状番号】 9722155

【プルーフの要否】 要

【書類名】 明細書

【発明の名称】 樹状図表示方法及び樹状図表示システム

【特許請求の範囲】

【請求項 1】 複数種類の生体高分子に対して複数の異なる条件で実験を行って得られたデータの組に基づいて前記複数の生体高分子のクラスタリング処理を行い、その結果を樹状図の形式で表示するステップと、

前記樹状図の部分木を選択するステップと、

選択された部分木を別ウィンドウで表示するステップとを含むことを特徴とする樹状図表示方法。

【請求項 2】 請求項 1 記載の樹状図表示方法において、

前記別ウィンドウに表示された部分木に含まれる生体高分子に対するクラスタリング手法の変更を指示するステップと、

指示されたクラスタリング手法によって前記部分木に含まれる生体高分子に対して再度クラスタリング処理を行い、その結果を樹状図の形式で表示するステップとを含むことを特徴とする樹状図表示方法。

【請求項 3】 複数種類の生体高分子に対して複数の異なる条件で実験を行って得られたデータの組に基づいて前記複数の生体高分子のクラスタリング処理を行い、その結果を樹状図の形式で表示するステップと、

前記樹状図の部分木を選択するステップと、

選択された部分木をアイコン化して表示するステップを含むことを特徴とする樹状図表示方法。

【請求項 4】 請求項 3 記載の樹状図表示方法において、アイコン化されて表示されている部分木を元の樹状図の形式に戻して再表示するステップを含むことを特徴とする樹状図表示方法。

【請求項 5】 複数種類の生体高分子に対して複数の異なる条件で実験を行って得られたデータの組に基づいて前記複数の生体高分子のクラスタリング処理を行い、その結果を樹状図の形式で表示するステップと、

前記樹状図の部分木を選択するステップと、

選択された部分木に含まれる生体高分子を対象として、生体高分子に関する情

報の中に予め用意されたキーワード辞書ファイルに格納されたキーワードが含まれている生体高分子の数を計数して表示するステップとを含むことを特徴とする樹状図表示方法。

【請求項 6】 複数種類の生体高分子に対して複数の異なる条件で実験を行って得られたデータの組に基づいて前記複数の生体高分子のクラスタリング処理を行い、その結果を樹状図の形式で表示するステップと、

前記樹状図の部分木を選択するステップと、

キーワードを指定するステップと、

生体高分子に関する情報の中に指定されたキーワードが含まれている生体高分子の前記部分木内での位置を表示するステップとを含むことを特徴とする樹状図表示方法。

【請求項 7】 請求項 1～6 のいずれか 1 項記載の樹状図表示方法において、前記生体高分子は cDNA、RNA、DNA 断片又は遺伝子であることを特徴とする樹状図表示方法。

【請求項 8】 複数種類の生体高分子に対して複数の異なる条件で実験を行って得られたデータの組に基づいて前記複数の生体高分子のクラスタリング処理を行い、その結果を樹状図の形式で表示するための解析を行うクラスタリング処理部と、

樹状図を表示するための表示部と、

入力手段と、

前記生体高分子に関する情報のキーワードを保持しているキーワード辞書ファイルとを備えることを特徴とする樹状図表示システム。

【請求項 9】 請求項 8 記載の樹状図表示システムにおいて、前記入力手段によって選択された部分木を別ウィンドウで表示する機能を有することを特徴とする樹状図表示システム。

【請求項 10】 請求項 9 記載の樹状図表示システムにおいて、前記別ウィンドウに表示された部分木に対してクラスタリング手法を変更して再度クラスタリング処理を行い、再クラスタリング処理によって得られた樹状図を表示する機能を有することを特徴とする樹状図表示システム。

【請求項 1 1】 請求項 8, 9 又は 1 0 記載の樹状図表示システムにおいて、前記入力手段によって選択された部分木をアイコン化して表示する機能、及びアイコン化されて表示されている部分木を元の樹状図の形式に戻して再表示する機能を有することを特徴とする樹状図表示システム。

【請求項 1 2】 請求項 8 ~ 1 1 のいずれか 1 項記載の樹状図表示システムにおいて、前記入力手段によって選択された部分木に含まれる生体高分子に対して、当該生体高分子に関する情報の中に前記キーワード辞書ファイルに格納されたキーワードが含まれている生体高分子の数を計数して表示する機能及び／又は選択されたキーワードを有する生体高分子の樹状図上の位置を表示する機能を有することを特徴とする樹状図表示システム。

【請求項 1 3】 請求項 8 ~ 1 2 のいずれか 1 項記載の樹状図表示システムにおいて、前記生体高分子は c D N A、R N A、D N A 断片又は遺伝子であることを特徴とする樹状図表示システム。

【発明の詳細な説明】

【0 0 0 1】

【発明の属する技術分野】

本発明は、特定の生体高分子、例えば遺伝子とハイブリダイズさせることによって得られたデータ（遺伝子発現データ）を、視覚的にわかりやすく、そして生体高分子（遺伝子）の機能・役割が推測しやすい形式によって表示するための表示方法及び表示システムに関する。

【0 0 0 2】

【従来の技術】

ゲノム配列が決定された種の増加に伴い、進化に対応すると見られる遺伝子を見つけ出し、どの生物にも共通に持っていると考えられる遺伝子の集合を探したり、それから逆に種に個別な特徴を推測するなど、種間の違いから何かを見出すとする、いわゆるゲノム比較法が盛んに行われてきた。しかし近年、DNAチップやDNAマイクロアレイなどのインフラストラクチャの発達によって、分子生物学の興味は、種間の情報から種内の情報へ、すなわち同時発生解析へと移りつつあり、これまでの種内の比較と併せて、情報の抽出から関連付けの場が大きく広

がりを持ち始めている。

【0003】

例えば、既知の遺伝子と同一の発現パターンを示す未知の遺伝子が見つければ、それが既知の遺伝子と同様の機能があると推測できる。これら遺伝子や蛋白質そのものの機能的な意味付けは、機能ユニットや機能グループといった形で研究されている。またそれらの間の相互作用も、既知の酵素反応データや物質代謝データとの対応づけによって、あるいはより直接的に、ある遺伝子を破壊あるいは過剰反応させ、その遺伝子の発現をなくすか、あるいは多量に発現させ、その遺伝子の直接的及び間接的影響を、全遺伝子の発現パターンを調べることによって解析している。

【0004】

この分野に成功した事例として、スタンフォード大学のP. Brownらのグループによるイースト菌の発現解析が挙げられる (Michel B. Eisen et al.: Cluster analysis and display of genome-wide expression patterns: Proc. Natl. Acad. Sci. (1998) Dec 8;95(25):14863-8)。彼らは、DNAマイクロアレイを用いて、細胞から抽出した遺伝子を時系列にハイブリダイズさせ、遺伝子の発現の度合い (ハイブリダイズした蛍光シグナルの輝度) を数値化した。そしてこの数値に応じて、細胞の一連のサイクルで発現パターンの過程が近い遺伝子どうし (任意の時点での発現の度合いが近いものどうし) をクラスタリングしている。

【0005】

図1は、この方式にそって遺伝子の発現パターンの類似性を表現した表示例である。右側には観測した個々の遺伝子の情報が列挙されており、左側にはこれらの遺伝子の発現パターンに応じて作成された樹状図が示されている。樹状図は、クラスタリングの過程で、最も近い2つのクラスタ毎に併合されてきた状況を表しており、各枝の長さは併合時の2つのクラスタ間距離 (クラスタ間の非類似度) に対応している。このような表示方法をとることで、共通のクラスタに属する遺伝子は、共通の機能的性質をもつ可能性があると推測することができる。

【0006】

【発明が解決しようとする課題】

実際の遺伝子発現パターンの分析では、大量のデータをクラスタリングすることになる。通常、DNAチップやDNAマイクロアレイは、数千から数万の遺伝子を同時に観測することが可能である。一般に遺伝子の発現過程は、ある遺伝子の発現が別の遺伝子の発現を誘導したり、あるいは、発現を阻害するなど、遺伝子間で複雑なネットワークを形成している。それ故、観測する遺伝子の数が多ければ、より複雑で詳細なネットワークを調べることができる。

【0007】

ところが、遺伝子の数が膨大になると、全体の遺伝子の働きを把握することは非常に困難になる。すなわち、樹状図には数千～数万の遺伝子が並ぶことになるので、この表示から、どのような分類ができているのかを判断するのは難しい。また、クラスタリング手法の違いにより、樹状図の枝の長さは一般的に異なる。例えばクラスタ併合アルゴリズムとして、最長距離法を選択したとき、枝の長さの平均は、最短距離法を選択したときの長さの平均よりも長い。したがって、樹状図全体としてみたとき、図2のように、根から葉までの長さもまた、クラスタリング手法によって異なる。遺伝子の発現データに対するクラスタリングでは、枝の長さよりも、どのように分類されているのかを調べることが重要である。それゆえ、通常、樹状図の表示を行なうときは、図3のように、樹状図の根から葉までの長さを一定値に定め、各枝の長さは根から葉までの長さに対する相対的な長さで表し、クラスタリング手法に応じて枝の長さの縮尺を変更して表示する。

【0008】

ここで、上記のような樹状図の表示方法を採用したとき、樹状図の中に発現パターンが類似している遺伝子を多数含んでいると、枝の長さが小さい樹状図が形成されるが、これらの枝の長さが樹状図の根から葉までの長さに対して非常に小さいと、図4の401に示すように遺伝子間の枝の詳細な関係を知るのが非常に困難になる。また、従来の遺伝子発現解析に関するクラスタリングでは、部分木を選択し、これに対して別のクラスタリング手法を適用するなど、対話的な操作ができなかった。また、従来の遺伝子発現解析に関するクラスタリングでは、分類が成功しているかどうかを調べる手段として、遺伝子の機能や遺伝子名のキーワードに着目し、それらが部分木に集まっているかどうかによって判断していた

。しかし、解析する遺伝子の数が膨大なものになると、どのような機能やキーワードに着目すべきかを判断するのは、非常に困難な作業である。

本発明は、このような従来技術の問題点に鑑み、樹状図全体の枝の状態を大域的に把握でき、かつ個々の部分木の状態を詳細に知ることができるような樹状図表示方法及び樹状図表示システムを提供することを目的とする。

【0009】

【課題を解決するための手段】

上記目的を達成するために、本発明では、樹状図の枝を選択し、選択した枝から葉の部分木に対して、別の表示ウィンドウで表示する機能、アイコン化する機能、アイコン化したものを元に戻す機能、部分木に含まれるキーワードを収集し表示する機能、を備えた樹状図表示システムを提案する。本発明によると、作成された樹状図の部分木に対して、異なるクラスタリング方法に対話的に適用する処理を実現することができる。また、クラスタリングが成功しているかどうかを判別するため、部分木にどのようなキーワードが多く含まれているかを表示し、分類の絞り込みや、クラスタリング方法の選択の支援を行うことができる。

【0010】

以下、理解を容易にするため、本発明を遺伝子のクラスタリングに適用した場合を例にとって、本発明の樹状図表示システムによる樹状図の表示例について説明する。ただし、本発明は遺伝子のクラスタリングにのみ適用されるものでなく、他の生体高分子、例えばcDNA、RNA、DNA断片等についても同様に適用可能である。

【0011】

図5は、本発明の樹状図表示システムによる樹状図の表示例を示している。分類アルゴリズムの選択メニュー501、及び（非）類似度の選択メニュー502を備えている。遺伝子発現データを読み込み、分類アルゴリズム及び（非）類似度を選択すると、樹状図が作成される。また、本システムは、図1のように遺伝子名などの遺伝子情報を樹状図の葉の先に付加して表示する形式も選択できる。

【0012】

作成された樹状図において、任意の枝を選択すると、選択した枝から葉までの

部分木に対する操作、すなわち、この部分木を別のウィンドウで表示する、この部分木をアイコン化する、この部分木のアイコンを元に戻す、この部分木に含まれる単語を検索する、というメニューが選択できる。図は、画面中央の枝 5 0 5 を矢印で図示されているマウスカーソル 5 0 4 等で選択した状態を示しており、このとき開くメニューウィンドウ 5 0 3 には選択可能なメニューが表示されている。マウスカーソル 5 0 4 をメニューウィンドウ 5 0 3 内に移動して、所望の処理項目をクリックすると選択された処理が実行される。

【0 0 1 3】

分類アルゴリズムは、図 5 の状態ではワード法が選択されているが、選択メニュー 5 0 1 を開くことによって例えば、最短距離法、最長距離法、群平均法、重心法、メディアン法、可変法など他のアルゴリズムを選択することができる。

(非)類似度は、個体間の類似の程度を表す指標である。この指標には、距離のように値の小さい方が類似性が高いことを表す場合と、相関係数のように値の大きい方が類似性が高いことを表す場合がある。前者の指標を非類似度、後者の指標を類似度という。図 5 の状態では非類似度としてユークリッド距離が選択されているが、選択メニュー 5 0 2 から他の(非)類似度、例えば標準化ユークリッド平方距離、マハラノビスの(汎)距離、ミンコフスキー距離等を選択することができる。このとき、分類アルゴリズムとして重心法、メディアン法、可変法を選択したとき、非類似度としてユークリッド平方距離以外に選択できないなど、分類アルゴリズムと非類似度との組み合わせが妥当なものである必要がある。

【0 0 1 4】

図 6 は、図 5 に示した表示画面において、「部分木を別のウィンドウで表示する」メニューを選択したときの表示例である。図 6 に示すように、選択した部分木を、根から葉までの長さに応じて縮尺を変更し表示し直す。このような表示手法をとることで、利用者は部分木の詳細な枝の状態を調べることが出来る。また、本システムでは、選択した部分木に対して、分類アルゴリズム及び／又は(非)類似度を選択して、再度クラスタリングを行なうことが出来る。このようにすることで、例えば、はじめのクラスタリング結果からクラスタ間の距離が大きいもの(図 4 において、4 0 1 と 4 0 2、4 0 1 と 4 0 3 の関係)を見つけ出し、

これを除外して、興味のある部分木のみ詳しく調べることが出来る。分類アルゴリズム及び／又は（非）類似度の選択は、分類アルゴリズムの選択メニュー 5 0 1、及び（非）類似度の選択メニュー 5 0 2 によって行う。

【0 0 1 5】

図 7 は、図 5 に示した表示画面において、「部分木をアイコン化する」メニューを選択したときの表示例である。部分木 5 0 5 を 7 0 1 のようにアイコンにすることで、樹状図の大域的な状態を容易に知ることが出来る。例えば、同様の機能をもつ遺伝子群や、発現がほとんど観測されなかった遺伝子群を一つのアイコンとしてまとめるなどの利用法が可能である。

【0 0 1 6】

図 8 は、図 5 に示した表示画面において、「部分木に含まれる単語を検索する」メニューを選択した時の表示例である。この機能を適用すると、選択した部分木に含まれる遺伝子の中で、遺伝子に対応する遺伝子情報に予め定めたキーワードが含まれるものを数え上げ、検索結果 8 0 1 として表示する。更に検索結果 8 0 1 から、マウスカーソル 8 0 4 等で一つのキーワード 8 0 2 を選択すると、そのキーワード（図の場合、"ribosomal"）を持つ遺伝子を、マーク 8 0 3 等によって樹状図上の位置として表示する。これにより、選択した部分木にどのような遺伝子が集まっているかを容易に知ることができる。また、この結果、分類がうまくいっていないのであれば、別の分類アルゴリズムや（非）類似度を選択して再度クラスタリングを行なうなど、より適切なクラスタリング方法の選択の支援をすることができる。

このように、本発明によると、作成された樹状図から、効果的に意味を抽出することができる。

【0 0 1 7】

すなわち、本発明による樹状図表示方法は、複数種類の生体高分子に対して複数の異なる条件で実験を行って得られたデータの組に基づいて前記複数の生体高分子のクラスタリング処理を行い、その結果を樹状図の形式で表示するステップと、前記樹状図の部分木を選択するステップと、選択された部分木を別ウィンドウで表示するステップとを含むことを特徴とする。

本発明は、別ウィンドウに表示された部分木に含まれる生体高分子に対するクラスタリング手法の変更を指示するステップと、指示されたクラスタリング手法によって前記部分木に含まれる生体高分子に対して再度クラスタリング処理を行い、その結果を樹状図の形式で表示するステップとを含んでもよい。

【0018】

本発明の樹状図表示方法は、また、複数種類の生体高分子に対して複数の異なる条件で実験を行って得られたデータの組に基づいて前記複数の生体高分子のクラスタリング処理を行い、その結果を樹状図の形式で表示するステップと、前記樹状図の部分木を選択するステップと、選択された部分木をアイコン化して表示するステップを含むことを特徴とする。

必要により、アイコン化されて表示されている部分木を元の樹状図の形式に戻して再表示するステップを含むこともできる。

【0019】

本発明による樹状図表示方法は、また、複数種類の生体高分子に対して複数の異なる条件で実験を行って得られたデータの組に基づいて前記複数の生体高分子のクラスタリング処理を行い、その結果を樹状図の形式で表示するステップと、前記樹状図の部分木を選択するステップと、選択された部分木に含まれる生体高分子を対象として、生体高分子に関する情報の中に予め用意されたキーワード辞書ファイルに格納されたキーワードが含まれている生体高分子の数を計数して表示するステップとを含むことを特徴とする。

【0020】

本発明による樹状図表示方法は、また、複数種類の生体高分子に対して複数の異なる条件で実験を行って得られたデータの組に基づいて前記複数の生体高分子のクラスタリング処理を行い、その結果を樹状図の形式で表示するステップと、前記樹状図の部分木を選択するステップと、キーワードを指定するステップと、

生体高分子に関する情報の中に指定されたキーワードが含まれている生体高分子の前記部分木内での位置を表示するステップとを含むことを特徴とする。

上記樹状図表示システムにおいて、生体高分子は cDNA、RNA、DNA 断片又は遺伝子とすることができる。

【0021】

本発明による樹状図表示システムは、複数種類の生体高分子に対して複数の異なる条件で実験を行って得られたデータの組に基づいて前記複数の生体高分子のクラスタリング処理を行い、その結果を樹状図の形式で表示するための解析を行うクラスタリング処理部と、樹状図を表示するための表示部と、入力手段と、生体高分子に関する情報のキーワードを保持しているキーワード辞書ファイルとを備えることを特徴とする。入力手段は、樹状図の枝の選択や、クラスタリング手法の選択などに用いられるもので、例えばキーボードやマウスとすることができる。キーワード辞書ファイルは、クラスタリングの結果に対し利用者が望む形になっているかを判断するために用いることができる。

【0022】

この樹状図表示システムは、入力手段によって選択された部分木を別ウィンドウで表示する機能を有することができる。また、別ウィンドウに表示された部分木に対してクラスタリング手法を変更して再度クラスタリング処理を行い、再クラスタリング処理によって得られた樹状図を表示する機能を有することができる。

この樹状図表示システムは、入力手段によって選択された部分木をアイコン化して表示する機能、及びアイコン化されて表示されている部分木を元の樹状図の形式に戻して再表示する機能を有することができる。

【0023】

この樹状図表示システムは、入力手段によって選択された部分木に含まれる生体高分子に対して、当該生体高分子に関する情報の中にキーワード辞書ファイルに格納されたキーワードが含まれている生体高分子の数を計数して表示する機能及び／又は選択されたキーワードを有する生体高分子の樹状図上の位置を表示する機能を有することができる。

本発明の樹状図表示システムにおいて、前記生体高分子はcDNA、RNA、DNA断片又は遺伝子とすることができる。

【0024】

【発明の実施の形態】

以下、図面を参照して本発明の実施の形態を説明する。以下では、遺伝子のクラスタリングを例にとって説明するが、本発明の適用範囲は遺伝子のクラスタリングのみに限定されるわけではなく、cDNA、RNA、DNA断片など生体高分子一般に対して同様に適用することができる。

【0025】

図9は、本発明による樹状図表示システムの一例を示す構成図である。このシステムは、遺伝子の情報及び発現過程を記録した遺伝子データ901と、遺伝子の発現過程に応じてクラスタリングを行ない、それを樹状図の形式で表示するための解析を行なうクラスタリング処理部902と、樹状図を表示するための表示装置903と、樹状図の枝や、クラスタリング手法の選択などに用いるキーボード904及びマウス905等の入力手段と、クラスタリングの結果に対し利用者が望む形になっているかを判断するための遺伝子情報のキーワードを保持しているキーワード辞書ファイル906から構成される。このクラスタリング処理部902は、コンピュータとそのプログラムによって具体化されるものである。なお、記憶装置901に代えて、ネットワーク等を介して遠隔地に設置されたサーバコンピュータが管理しているデータベースから遺伝子データを取得する構成をとってもよい。

【0026】

図10は、遺伝子データ901に格納された遺伝子発現パターンデータの具体的な構造を示したものである。本アルゴリズムでは、これを2次元配列によって格納する。すなわち、遺伝子ID(id)をもつ遺伝子が実験ケース(no)における発現の度合い(ハイブリダイズした蛍光シグナルの輝度)を数値化したデータを、Exp[id][no]に格納する。m種類の遺伝子をそれぞれ異なる位置にスポットしたDNAチップから得られる1回の実験は、1つの実験ケースに対応する。

【0027】

図11は、遺伝子データ901に格納された遺伝子に関する情報を格納するための、遺伝子情報構造体の例を示している。この遺伝子構造体は、遺伝子ID(1101)、遺伝子のORF(1102)、遺伝子名(1103)、遺伝子の機能(1104)のメンバから構成される。図11はあくまでも説明のための例であり

、ここに示した遺伝子の属性以外の情報も、遺伝子情報構造体のメンバとして定義することももちろん可能である。

【0028】

図12は、クラスタリング処理において利用するクラスタを表す構造体の例を示している。全てのクラスタ構造体は、樹状図の各ノードまたは葉と対応している。クラスタ構造体は、ウィンドウ単位で管理され、同じウィンドウのノードまたは葉であれば、同一のwindowID (1207) をもつ。また、同じウィンドウ内でノードまたは葉を識別するため、clusterNo (1205) で各クラスタ構造体に一意に番号を割り振っている。クラスタ構造体には3種類あり、type (1201) の値が、leafのもの、nodeのもの及びiconのものに分かれる。

【0029】

各leaf型クラスタ構造体は、それぞれひとつのgeneID (1206) に対応している。すなわち、ひとつの遺伝子に対応している。更にgeneIDから、遺伝子情報構造体のデータが参照できる。node型クラスタ構造体は、クラスタリングにおける併合処理において逐次生成するもので、併合前の2つのクラスタをleft (1202) の値と、right (1203) の値からたどれるようにし、また、それらの間の距離 ((非) 類似度) をdistance (1204) の値として保持する。left及びrightの値には、クラスタを一意に示すclusterNo (1205) が入っている。icon型クラスタ構造体は、部分木をアイコンに置き換えるときに生成され、表示では葉の場合と同様に扱う。そして枝の先端には部分木を示すアイコンを付して表示する。実際の部分木のルートのクラスタは、left (1202) の値からたどることができる。

【0030】

図13は、図12に例示したクラスタ構造体のデータ構造を示した図である。

これらはクラスタ分析の過程で生成される。クラスタ構造体は、最初leaf型のものだけを用意するが、クラスタリングの過程で2つずつ併合し、その度にnode型クラスタ構造体を生成してトリートメント構造を組み立てる。node型クラスタ構造体には、併合した2つの子ノードのclusterNoと、それらの間の距離 ((非) 類似度) の情報が登録されている。またleaf型クラスタ構造体に登録されているgeneIDによ

り、対応する遺伝子情報を参照することができる。アイコン化する処理があれば、トリーの途中にicon型のクラスタを挿入し、あたかも葉であるかのように表示する（表示に当たっては、icon型のクラスタより先に位置するクラスタは表示しない）。アイコンを解除するときは、icon型クラスタの上下のクラスタのリンクを繋ぎ直す操作を行う。

【 0 0 3 1 】

図 1 4 は、クラスタ分析の過程でクラスタ間の距離である非類似度を格納するための配列の例である。図に示すように、2次元配列dist[][]を用いてこれを格納する。また、2次元配列dist[][]のインデックスの数字に対応するクラスタのclusterNo (1 2 0 5) を格納した配列を、clust_idx[]に格納する。すなわち、非類似度dist[i][j]の値は、clusterNoがclust_idx[i]とclust_idx[j]であるクラスタ間の値を示す。図 1 4 から、例えばclust_idx[3]であるclusterNo: 9のクラスタとclust_idx[4]であるclusterNo: 25のクラスタ間の非類似度dist[3][4]の値は21であることが分かる。

【 0 0 3 2 】

図 1 5 は、各ウィンドウのルートノードを格納する配列の例を示している。すなわち、各表示ウィンドウに対するルートノードのクラスタのclusterNoは、配列RootNode[]に格納される。図 1 5 に示した例では、RootNode[1]の値が569であることからwindowID: 1の表示ウィンドウに表示される樹状図のルートノードはclusterNo: 569のクラスタであることが分かり、RootNode[2]の値が312であることからwindowID: 2の表示ウィンドウに表示される樹状図のルートノードはclusterNo: 312のクラスタであることが分かる。

【 0 0 3 3 】

図 1 6 は、検索の問合せ及び結果を格納するためのsearch構造体の例を示している。キーワード辞書ファイル906に登録されている各キーワードに対して、構造体の一つ生成する。また、キーワードで同義語のものがいくつか存在するとき、それらをひとつのものを指すこととして扱うこともできる。search構造体は、検索項目のキーワードを入力しておくkeyword (1 6 0 1)、そのキーワードが部分木の中でいくつあったかを示すtimes (1 6 0 2)、キーワードが遺伝子

情報の中にあったとき、その遺伝子の樹状図上の位置を格納するplace (1 6 0 3) をメンバとしてもつ。図 1 6 に図示する例のように、Rat、Mouse、Musのような同義語をまとめてkeywordメンバに登録しておくことで、これら3つのキーワードのどれをも同じ検索項目として扱うことが出来る。

【0 0 3 4】

図 1 7 は、本システムの概略フローを示した図である。

まず、遺伝子データ 9 0 1 からクラスタリング処理部 9 0 2 へデータを読み込む (ステップ 1 7 0 1)。これについては、後で詳しく説明する。次に、クラスタ分析、及び結果表示に必要な各種パラメータを設定する (ステップ 1 7 0 2)。ここでは、分類アルゴリズム及び (非) 類似度の設定、個々の遺伝子情報を表示するか否かなどの設定を行う。

【0 0 3 5】

次にクラスタ分析を行い (ステップ 1 7 0 3)、結果を表示する (ステップ 1 7 0 4)。クラスタ分析については、後で詳しく説明する。このクラスタ分析の処理の中で、樹状図表示に必要な情報を収集し、クラスタ構造体に入力する。分析結果表示では、このクラスタ構造体と、個々のウィンドウのルートノードのclusterNoを表すRootNode[]の情報をもとに、結果を表示する。クラスタ構造体のtypeがiconのときは、それを葉のように扱い、枝の先端に部分木を表すアイコンを付加する。

【0 0 3 6】

表示された樹状図の中のある部分木をアイコン化してまとめる、あるいはアイコン化を解除して元の部分木に戻す場合、以下の処理を実行する (ステップ 1 7 0 5)。すなわち、樹状図の枝をマウスで選択し (ステップ 1 7 0 6)、部分木のアイコン化、または非アイコン化処理を行う (ステップ 1 7 0 7)。アイコン化、非アイコン化処理に関しては、後で詳しく説明する。処理の後、再び分析結果表示 (ステップ 1 7 0 4) を行う。

【0 0 3 7】

表示された樹状図に対して、キーワード辞書ファイル 9 0 6 に格納されたキーワードをもとに検索を行う場合、以下の処理を実行する (ステップ 1 7 0 8)。

すなわち、樹状図の枝をマウスで選択し（ステップ1709）、検索処理を行う（ステップ1710）。検索処理に関しては、後で詳しく説明する。検索処理1710で、表示に必要な情報がsearch構造体に格納されるので、それをもとに新たに検索結果ウィンドウを生成し結果を表示する（ステップ1711）。このとき、マウスなどで検索結果ウィンドウのあるキーワードを選択すると、search構造体のplaceメンバの情報をもとに、樹状図上のキーワードのある箇所にマーカーを付与する。

【0038】

表示された樹状図に対して、他の併合アルゴリズム、（非）類似度で再びクラスタリングを適用したいときは、ステップ1702に戻る（ステップ1712）。クラスタ併合アルゴリズムとしては、例えば、最短距離法、最長距離法、群平均法、重心法、メディアン法、ウォード法、可変法等がある。最短距離法、最長距離法、群平均法、ウォード法、可変法には、次々にクラスターを融合していくときの非類似度が単調に大きくなる特性がある。また、2つのクラスターを融合して1つのクラスターを作ると、他のクラスターとの距離が近づく場合と遠ざかる場合があり、前者を空間の収縮、後者を空間の膨張、距離が変わらない場合を空間の保存と呼ぶが、最短距離法は空間が収縮する特性を有し、最長距離法やウォード法は空間が膨張する特性を有する。また、群平均法、重心法、メディアン法は、空間が保存され、可変法の場合はパラメータの設定によっていずれにもなりうる。（非）類似度にも種々のものがあり、例えば非類似度の代表的なものとしてはユークリッド平方距離、標準化ユークリッド平方距離、マハラノビスの（汎）距離、ミンコフスキー距離等がある。従って、前述の特性等を勘案して、これらの中から適宜のものを選択すればよい。

【0039】

表示された樹状図に対して、ある部分木を別のウィンドウで表示させたい時（ステップ1713）は、別ウィンドウに表示したい樹状図の枝をマウスで選択し（ステップ1714）、選択した樹状図の部分木に対するデータの読み込みを行い（ステップ1715）、再びステップ1702に戻る。選択した樹状図の部分木に対するデータの読み込み処理については、あとで詳しく説明する。

以上の選択が無かった場合には、処理を終了する。

【0040】

図18は、図17における遺伝子データの読み込み処理1701の詳細フローである。

まず、遺伝子数、実験ケースの総数をそれぞれgene_num、exp_numに登録する（ステップ1801）。次に、遺伝子データ901から遺伝子情報を読み取り、遺伝子情報構造体gene_info[i] ($i = 1, \dots, \text{gene_num}$)に登録する（ステップ1802）。遺伝子データ901から遺伝子発現データを読み取り、Exp[i][j] ($i = 1, \dots, \text{gene_num}, j = 1, \dots, \text{exp_num}$)に登録する（ステップ1803）。樹状図の葉の総数を表すleaf_numにgene_numを代入する（ステップ1804）。

【0041】

次に、初期値となるleaf型クラスタ構造体を生成する。クラスタ構造体clusterをleaf_num個生成し、 $i = 1, \dots, \text{leaf_num}$ に対して、typeをleafに、clusterNoをiに、geneIDをiに、windowIDを1として登録する（ステップ1805）。次に、キーワード辞書ファイル906に格納されたキーワードを読み出し、それぞれのキーワードに対してsearch構造体を生成し、キーワードをsearch[].keywordに登録する（ステップ1806）。キーワードの総数をkey_numに代入する（ステップ1807）。windowIDを表すwidに1を登録し（ステップ1808）、処理を終わる。

【0042】

図19、図20は、図17におけるクラスタ分析処理1703の詳細フローである。

windowIDがwidに対応するウィンドウ内の遺伝子間の発現度の非類似度を求める。clusterNoがi,jに対応する遺伝子の非類似度をdist[i][j]に登録する（ステップ1901）。本アルゴリズムでは、クラスタが1つ生成されるごとにclusterNoを1から順に割り振っている。そこで、次のクラスタが生成されたとき、そのクラスタの番号を表すnewclusterNoにleaf_num + 1を代入しておく（ステップ1902）。また、クラスタ間距離（非類似度）を格納する配列の情報として、併合対象クラスタ数を示すall_clustにleaf_numを代入し、 $i = 1, \dots, \text{leaf_num}$ に

対し、`cluster_idx[i]` に i を代入して初期化しておく。併合対象クラスタの数 `all_clust` が 1 に等しいかどうか判定し、等しくない場合、1 になるまで以下の一連の処理を繰り返す（ステップ 1905）。

【0043】

最初に、先に求めたクラスタ間距離（非類似度）から、次に併合されるべきクラスタを決定する。すなわち、 $i < j$ かつ $i, j = 1, 2, \dots, \text{all_clust}$ に対して、`dist[i][j]` の最小値、最小値を与える i 、最小値を与える j を求め、`d_min`、`i_min`、`j_min` にそれぞれ代入する。`clusterNo` が、`cluster_idx[i_min]`、`cluster_idx[j_min]` のクラスタが次に併合されるべきクラスタとなる。`cluster` を新規に生成し、`type` に `node`、`left` に `cluster_idx[i_min]`、`right` に `cluster[j_min]`、`distance` に `d_min`、`clusterNo` に `newclusterNo`、`windowID` に `wid` を登録していく（ステップ 1907）。ここで、2 つのクラスタのどちらを `left` メンバとし、残りを `right` メンバとするかについては、発現量で比較するなど予め判定基準を設ける方式をとることも可能である。

【0044】

次に、クラスタ間距離を格納している配列の情報を更新する。まず、新しく生成したクラスタと他のクラスタとの距離（（非）類似度）を求め、それを `i_min` のクラスタと他のクラスタ間の距離が格納されていた `dist[][]` の配列位置に上書きする。 $i = 1, 2, \dots, i_{\min}-1$ に対し、新しく生成したクラスタと、`clusterNo` が `cluster_idx[i]` に対応するクラスタとの非類似度を `dist[i][i_min]` に登録し、 $j = i_{\min} + 1, \dots, j_{\min}-1, j_{\min} + 1, \dots, \text{all_clust}$ に対し、新しく生成したクラスタと、`cluster_idx[j]` に対応するクラスタとの非類似度を `dist[i_min][j]` に登録する（ステップ 2001、2002）。

【0045】

次に、`j_min` に関する情報を削除して、`j_min` 以降のすべての配列データを一つ前に移動する処理を行なう。 $i = \text{min_j}, \dots, \text{all_clust}-1$ に対し、`clust_idx[i]` に `clust_idx[i+1]` を代入する（ステップ 2003）。次に $i < j$ 、 $i, j = j_{\min}, \dots, \text{all_clust}$ を満たす i, j に対し、`dist[i][j]` に `dist[i+1][j]` を代入し、その後 $i < j$ 、 $i = 1, \dots, \text{all_clust}-1$ 、 $j = j_{\min}, \dots, \text{all_clust}-1$ を満たす

i, j に対し、 $\text{dist}[i][j]$ に $\text{dist}[i][j+1]$ を代入する（ステップ 2004、2005）。

【0046】

最後に、併合対象クラスタ数を示す all_clust から 1 を引き、新しいクラスタ構造体に割り振る clusterNo を表す newclusterNo に 1 を加える（ステップ 2006、2007）。

以上の操作を all_clust が 1 になるまで繰り返す。 all_clust が 1 になれば、 $\text{RootNode}[\text{wid}]$ に、このウィンドウのルートノードの clusterNo を表す $\text{cluster_idx}[1]$ を代入し、処理を終える（ステップ 1908）。

【0047】

図 21 は、図 17 におけるアイコン化する、または（非）アイコン化（アイコンを解除）する処理 1707 の詳細フローである。

6 において選択した枝の両端に対応するクラスタを登録する。下（leaf 側）の cluster を childClust に代入し、枝の上（root 側）の cluster を parentClust に代入する（ステップ 2101、2102）。次に、新しく icon 型 cluster を生成し、 childClust と parentClust の間に挿入する処理を行なう。すなわち、 cluster を生成し、 type に icon を、 left に $\text{childClust.clusterNo}$ 、 clusterNo に newclusterNo を、 windowID に wid をそれぞれ登録する（ステップ 2103）。そして、ポインタの付け替え操作として、 parentClust.left または parentClust.right に登録されている childClust の clusterNo を newclusterNo に変更する（ステップ 2104）。全体のクラスタ数がひとつ増加したので、新しいクラスタ構造体に割り振る clusterNo を示す newclusterNo に 1 を加えて処理を終了する。（ステップ 2105）

【0048】

また、部分木をアイコン化したものを元に戻すメニューを選択すると、まず図 17 におけるステップ 1706 で選択した枝の両端に対応するクラスタを登録する。ステップ 1706 で選択した枝の下（leaf 側）にあるアイコンの cluster 、アイコンの親ノードの cluster をそれぞれ iconClust 、 parentClust に代入する（ステップ 2101、2106）。アイコンのクラスタと、部分木のクラスタとの

ポインタを繋ぎ替え、アイコンのクラスタを削除する処理を行なう。すなわち、parentClust.leftまたはparentClust.rightに登録されているiconClustのclusterNoをiconClust.leftに変更する（ステップ2107）。その後、iconClustを削除して処理を終了する（ステップ2108）。

【0049】

図22は、図17における検索処理1710の詳細フローである。

選択した枝以下に対応する部分木のルートノードのクラスタのclusterNoをclustNoに代入する（ステップ2201）。また、部分木の先頭からのインデックスを表すleafNoを1で初期化しておく（ステップ2202）。また $i = 1, \dots, \text{key_num}$ に対して、search[i].timesを0、search[i].placeをnullで初期化しておく（ステップ2203）。次に、再帰的にクラスタ木に対するトリートワークを実行し、searchで指定したキーワードをもつ遺伝子の単語検索処理（処理A）を行なう（ステップ2205）。引数としてclustNo、leafNoを渡す。単語検索処理については、後で詳しく説明する。処理Aを終えると、search構造体に検索結果が入力され、処理を終了する。

【0050】

図23は、図22の単語検索処理（処理A）の詳細フローである。

引数で渡されたclustNo、leafNoをそれぞれclustNo、leafNoに代入する（ステップ2300）。また、clusterNoの指すclusterをtargetClustに代入する（ステップ2301）。キーワード検索のカウンタを示すiを0に設定しておく（ステップ2302）。

【0051】

次に、targetCluster.typeがleafかどうかを判定する（ステップ2303）。leafであるとき、leafに対応する遺伝子情報とキーワード辞書ファイルから読み込んだキーワードとの比較が終わるまで、以下の処理を繰り返し行なう。すなわち、iがkey_numになるまで繰り返し行なう（ステップ2304）。まず、targetClust.geneIDのgeneIDに対応する遺伝子情報構造体gene_infoの属性の中に、search[i].keywordの用語が入っているか判別する（ステップ2305）。もし入っていたら、部分木でキーワード（search[i].keyword）が発見された回数を示すs

earch[i].timesをひとつインクリメントし、部分木での発見した位置のインデックスを示すsearch[i].placeに現在位置のleafNoを登録する（ステップ2307）。キーワードの検索カウンタiをひとつインクリメントし、ステップ2304に戻る。ステップ2304において、iがkey_numになったとき、即ちすべてのキーワードとの比較が終わったら、部分木のインデックスであるleafNoをひとつインクリメントし、処理を終わる（ステップ2309）。

【0052】

また、ステップ2303において、targetCluster.typeがleafではなかった場合、子供のノードをたどる処理を行なう。targetClust.leftをclustNoに代入し（ステップ2310）、左の子ノードに対しclustNoとleafNoとを引数として再び単語検索処理（処理A）を行なう（ステップ2311）。targetCluster.typeがiconのときは、targetCluster.rightには子供ノードがないので、処理を終了する（ステップ2312）。ステップ2312において、targetCluster.typeがiconでない場合、これはnode型clusterを表す。clustNoにtargetClust.rightを代入し（ステップ2313）、右の子ノードに対しclustNoとleafNoとを引数として再び単語検索処理（処理A）を行ない、処理を終了する（ステップ2314）。

【0053】

図24は、図17における部分木の遺伝子データの読み込み処理1715の詳細フローである。

新しく部分木を読み込んでウィンドウを作成するので、新しいウィンドウIDを示すwidをひとつインクリメントしておく（ステップ2401）。また、樹状図の葉の総数を表すleaf_numを0に初期化しておく（ステップ2402）。選択した枝以下に対応する部分木のルートノードのクラスタにおけるclusterNoをclustNoに代入する（ステップ2403）。最後に、部分木のleaf型クラスタに対して、新規clusterを生成する処理（処理B）を行なう（ステップ2404）。現在のクラスタを示すclustNoをこの処理の引数として渡す。この処理の詳細は後で説明する。すべてのleafを読み込み、leafに対応するclusterをすべて生成し処理を終了する。

【0054】

図25は、図24における部分木のleafに対して新規にクラスタを生成する処理2404の詳細フローである。

引数で渡されたclustNoをclustNoとし、clustNoの指すclusterをtargetClustとする（ステップ2501、2502）。次に、targetCluster.typeがleafかどうかを判定する（ステップ2503）。leafであるならば、部分木のleafの数のカウンタであるleaf_numをひとつインクリメントする（ステップ2504）。次に新しいウィンドウの初期値となるleaf型クラスタ構造体を生成する。すなわち、clusterを生成し、typeにleafを、clusterNoにleaf_numを、geneIDにtargetCluster.geneIDを、windowIDにwidを登録し処理を終了する（ステップ2505）。

【0055】

またステップ2503において、targetCluster.typeがleafではなかった場合、子供のノードをたどる処理を行なう。すなわち、targetClust.leftをclustNoに代入し（ステップ2506）、左の子ノードに対し、clustNoを引数として再び新規にクラスタを生成する処理（処理B）を行なう（ステップ2507）。targetCluster.typeがiconのときは、targetCluster.rightに子供ノードはないので、これで処理を終了する（ステップ2508）。ステップ2508において、targetCluster.typeがiconでない場合、これはnode型clusterを表している。従って、clustNoにtargetClust.rightを代入し（ステップ2509）、右の子ノードに対しclustNoを引数として再び新規にクラスタを生成する処理（処理B）を行い、処理を終了する（ステップ2510）。

以上では解析結果を表示装置画面に表示する例を説明したが、多色プリンタで印刷出力する構成であってもよい。すなわち、本発明でいう表示とは、プリンタによって視覚的に印刷出力する概念を含むものである。

【0056】

【発明の効果】

以上示したように、本発明によると、樹状図に対して様々なクラスタリング手法を適用し、部分木をアイコン化したり、別ウィンドウで表示するなど、遺伝子

の発現解析等を支援する方法を提供することができる。

【図面の簡単な説明】

【図 1】

標準的クラスタ分析結果の表示例を示す図。

【図 2】

クラスタリング方法の違いの例の説明図。

【図 3】

クラスタリング方法によらない樹状図の表示例を示す図。

【図 4】

発現パターンが類似している遺伝子群を含む樹状図の例を示す図。

【図 5】

本発明の樹状図表示システムによる画面表示例を示す図。

【図 6】

本発明の樹状図表示システムによる他の画面表示例を示す図。

【図 7】

本発明の樹状図表示システムによる他の画面表示例を示す図。

【図 8】

本発明の樹状図表示システムによる他の画面表示例を示す図。

【図 9】

本発明による樹状図表示システムの構成例を示す図。

【図 10】

遺伝子発現パターンデータの例を示す図。

【図 11】

遺伝子情報構造体の例を示す図。

【図 12】

クラスタ構造体の例を示す図。

【図 13】

クラスタ木構造の生成例を示す図。

【図 14】

クラスタ間距離を格納する配列の例を示す図。

【図 1 5】

各ウィンドウのルートノードを格納する配列の例を示す図。

【図 1 6】

検索の問合せ及び結果を格納する構造体の例を示す図。

【図 1 7】

本システムの概略処理フロー例を示す図。

【図 1 8】

遺伝子データの読み込み処理のフローを示す図。

【図 1 9】

クラスタ分析処理のフローを示す図。

【図 2 0】

クラスタ分析処理のフローを示す図。

【図 2 1】

(非) アイコン化処理のフローを示す図。

【図 2 2】

遺伝子情報を検索対象とした検索処理のフローを示す図。

【図 2 3】

単語検索処理 (処理A) のフローを示す図。

【図 2 4】

部分木の遺伝子データの読み込み処理の説明図。

【図 2 5】

部分木の leaf に対して新規に cluster を生成する処理 (処理B) の説明図。

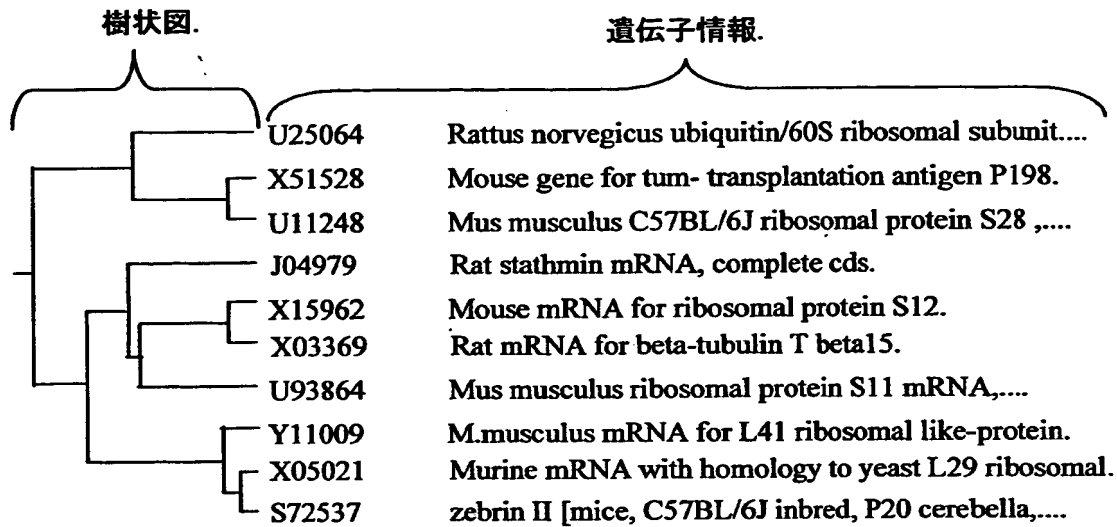
【符号の説明】

4 0 1 … 樹状図の中で発現過程が類似した遺伝子群の例、4 0 2 … 樹状図の中で発現過程が 4 0 1 の遺伝子群と大きく異なる遺伝子の例 (その 1)、4 0 3 … 樹状図の中で発現過程が 4 0 1 と大きく異なる遺伝子の例 (その 2)、5 0 1 … クラスタリングにおける分類アルゴリズムの選択メニュー、5 0 2 … クラスタリングにおける (非) 類似度の選択メニュー、5 0 3 … メニューウィンドウ、5 0

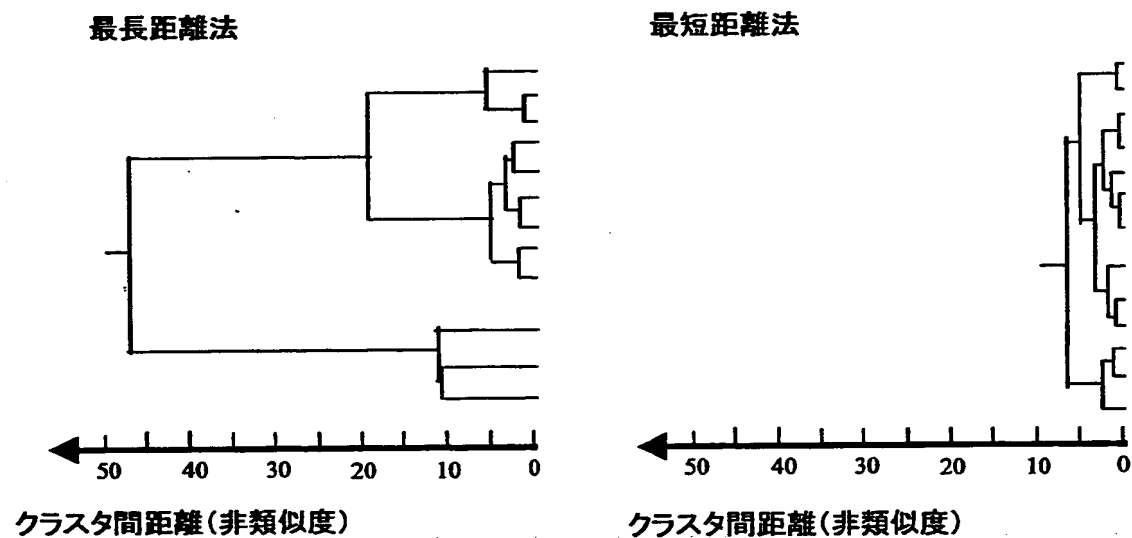
4…マウスカーソル、5 0 5…選択された枝（部分木）、7 0 1…アイコン化した部分木の例、8 0 1…キーワード検索結果のウィンドウ例、8 0 2…選択されたキーワード、8 0 3…遺伝子情報の中に予め定めたキーワードが含まれる遺伝子に対するマーク、8 0 4…マウスカーソル

【書類名】 図面

【図 1】

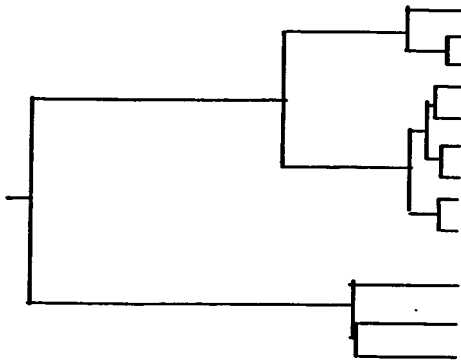


【図 2】

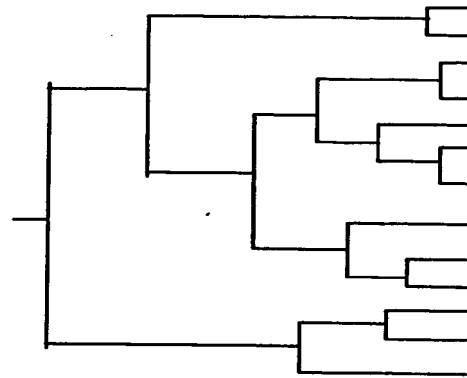


【図 3】

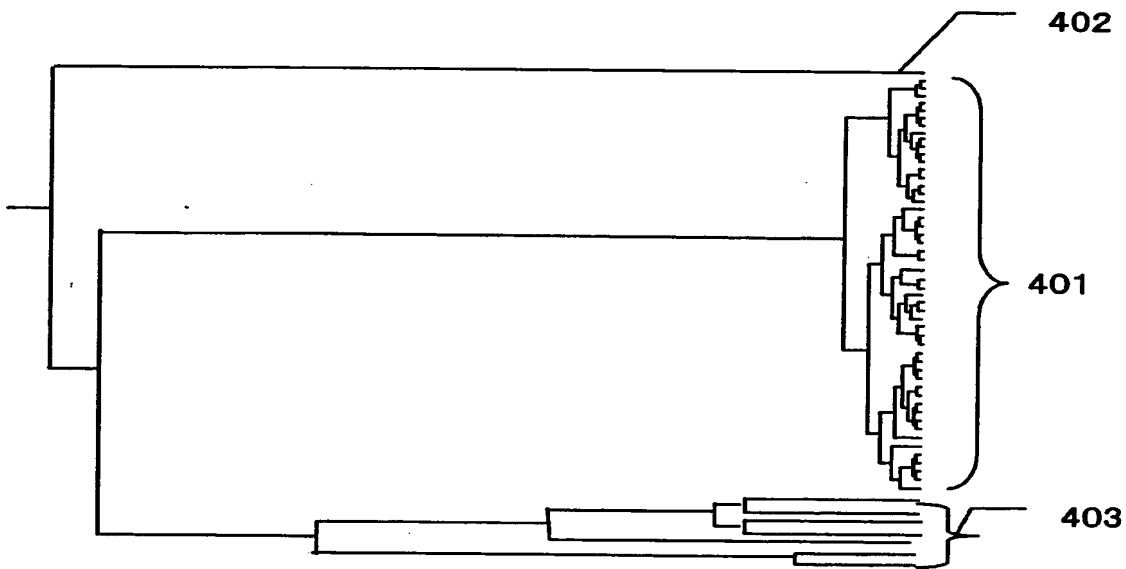
最長距離法



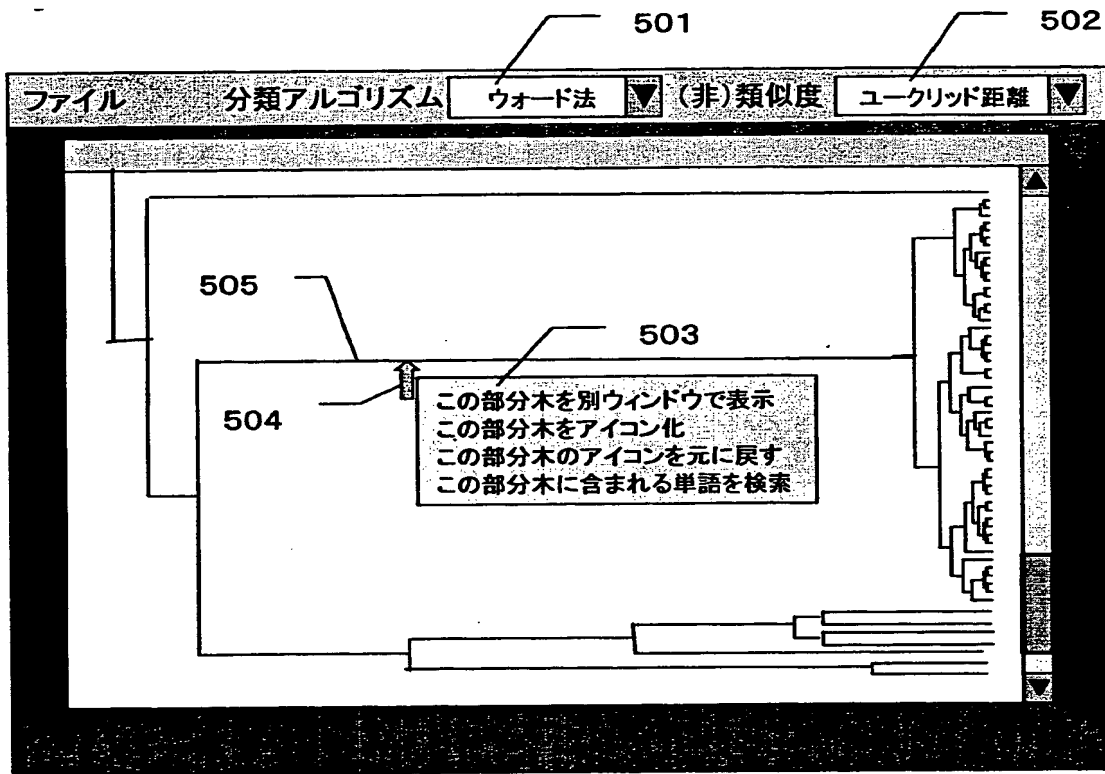
最短距離法



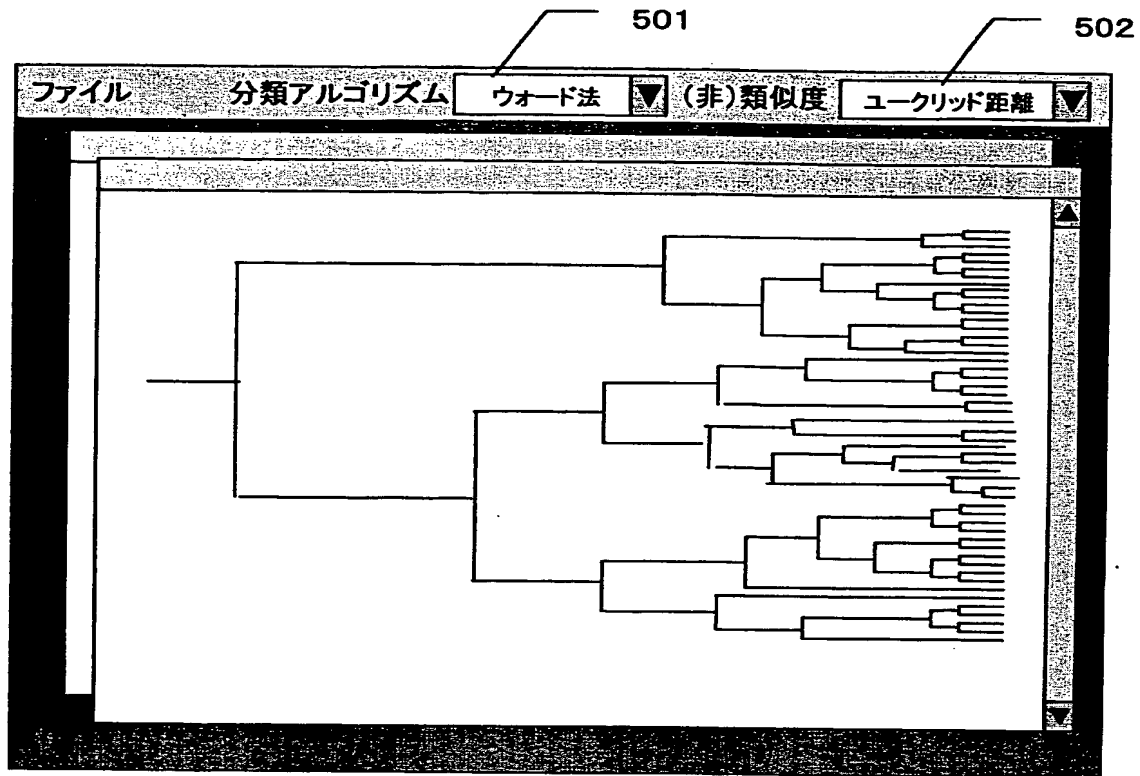
【図 4】



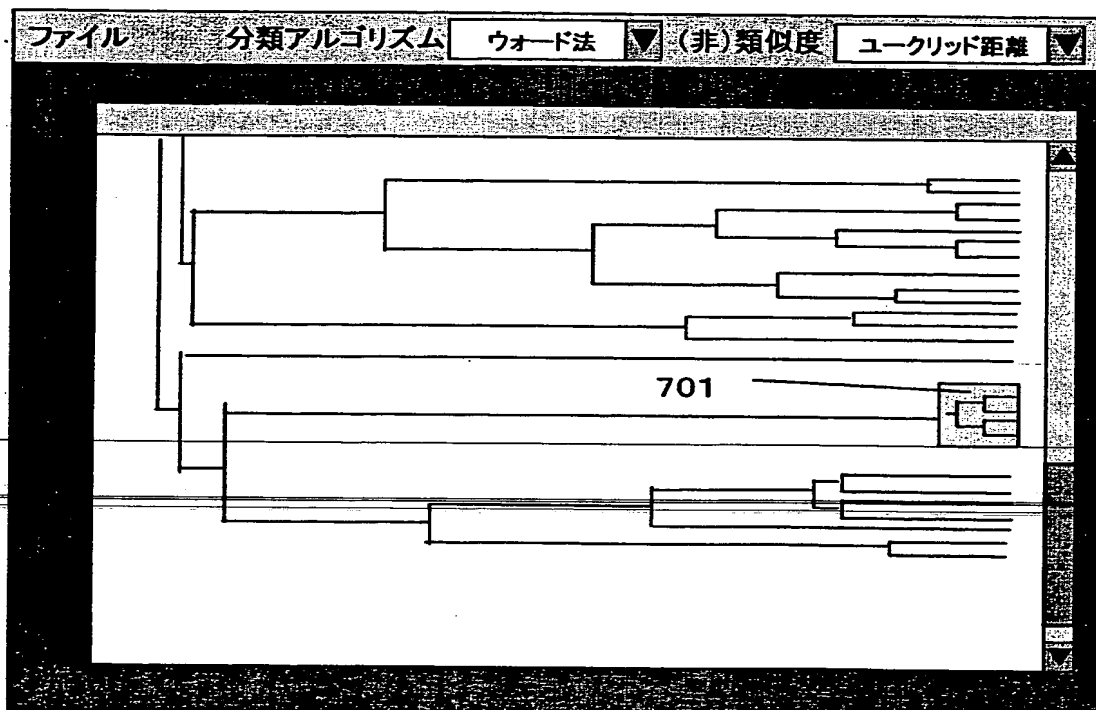
【図 5】



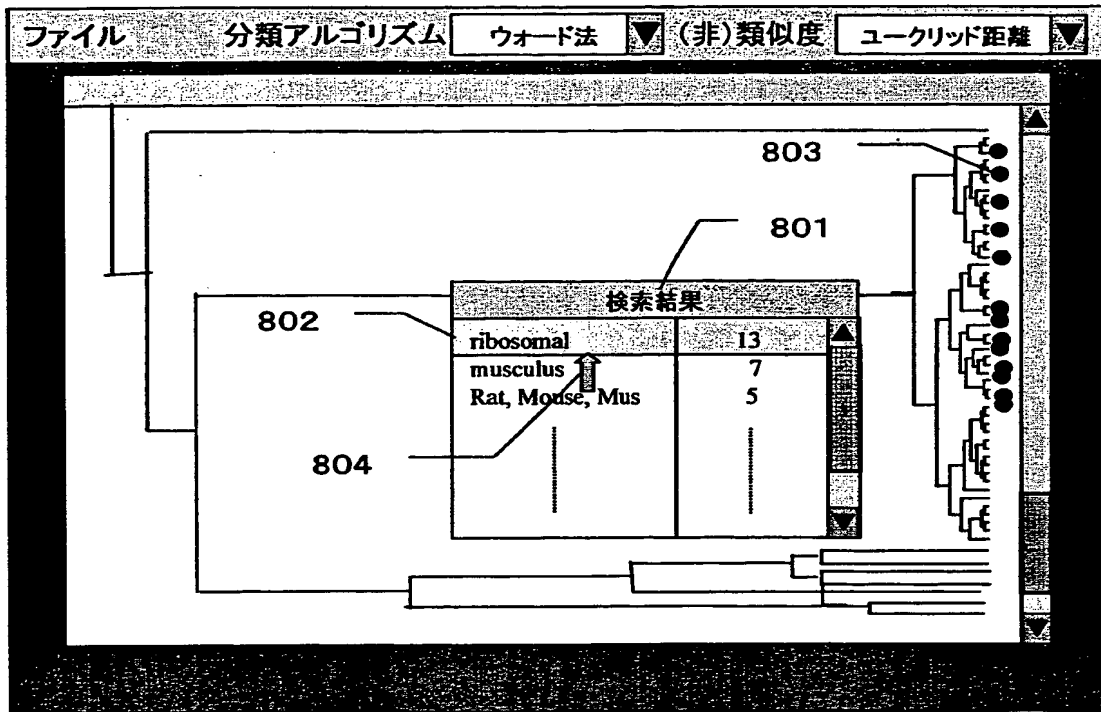
【図 6】



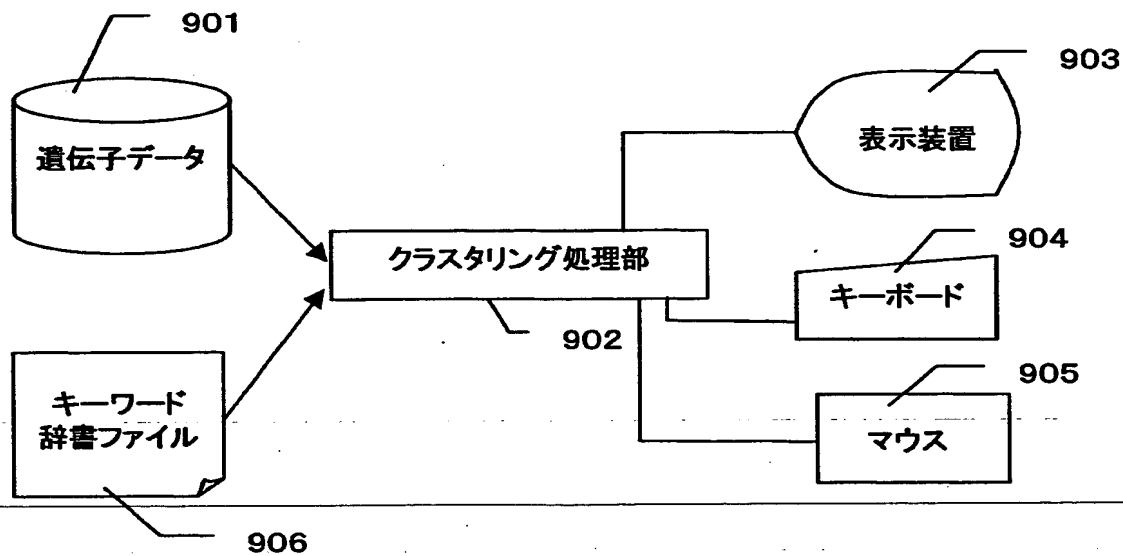
【図 7】



【図 8】



【図 9】



【図 1 0】

遺伝子ID (geneID)	実験ケース									
	1	2	3	4	5	6	..	no	..	n
1	0	1	2	0	3	2	0			
2	1	2	0	0	2	2	1			
:										
id	0	4	3	6	5	4	0			
:										
m	0	4	3	6	5	4	0			

Exp[id][no]

【図 1 1】

gene_info	
メンバ名	値
1101 geneID	17
1102 ORF	YBL084C
1103 name	ANAPHASE-PROMOTING COMPLEX SUBUNIT
1104 function	CELL CYCLE

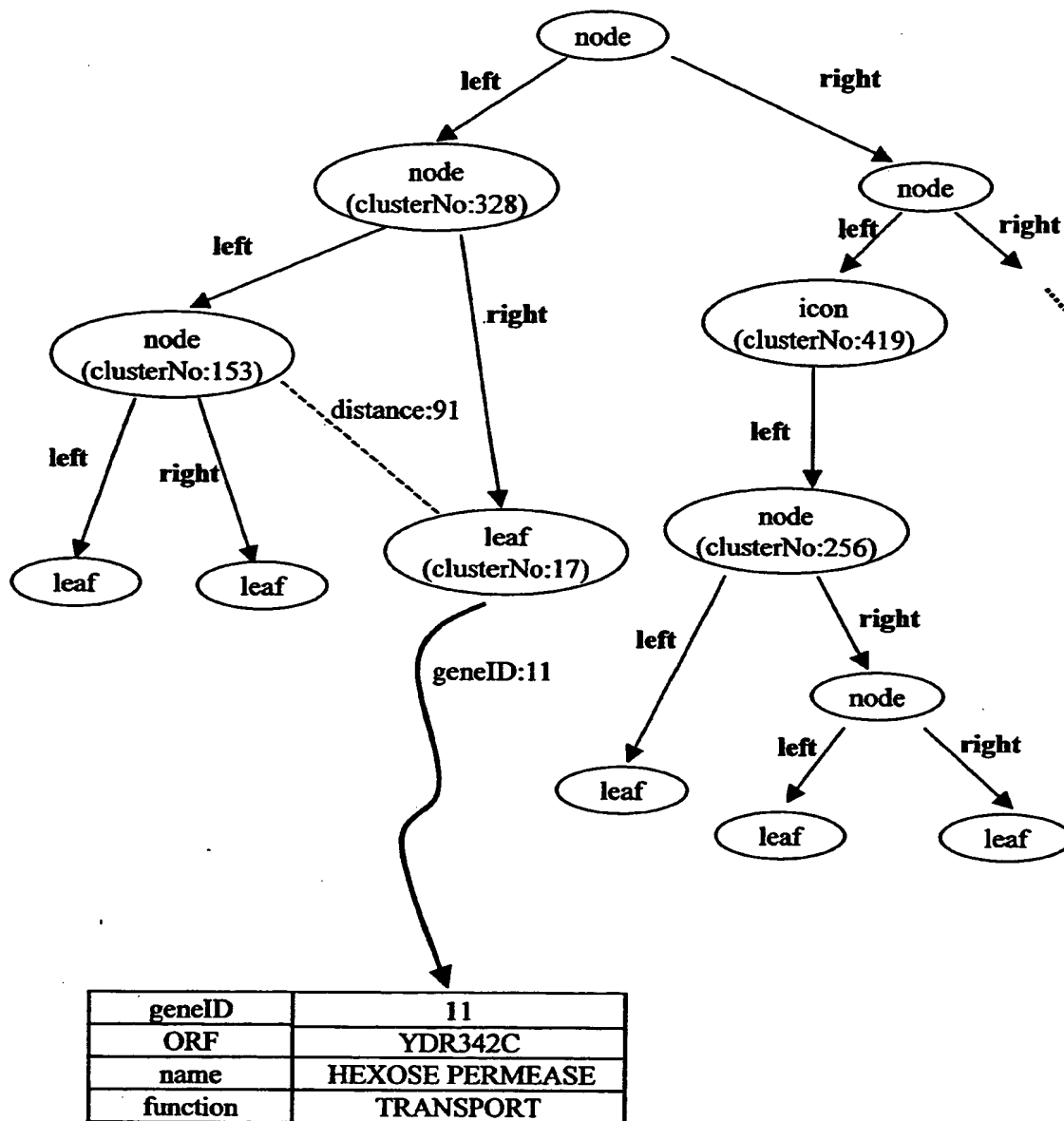
【図 1 2】

cluster	
	メンバ名
	値
1201	type
1202	left
1203	right
1204	distance
1205	clusterNo
1206	geneID
1207	windowID
	leaf
	17
	11
	3

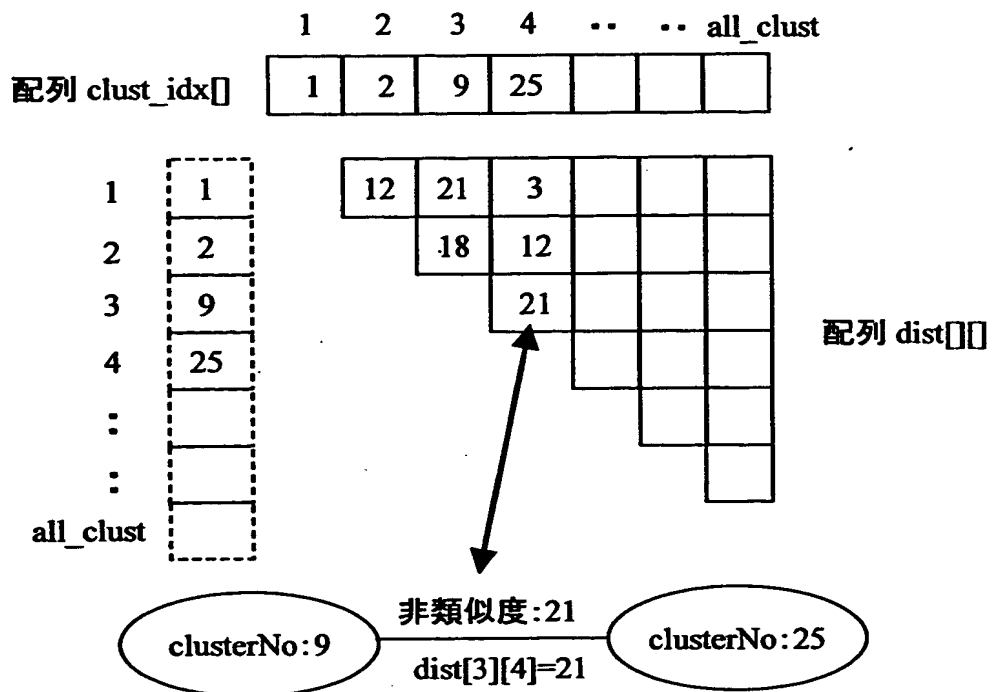
cluster	
	メンバ名
	値
	type
	left
	right
	distance
	clusterNo
	geneID
	windowID
	node
	153
	17
	91
	328
	3

cluster	
	メンバ名
	値
	type
	left
	right
	distance
	clusterNo
	geneID
	windowID
	icon
	256
	419
	3

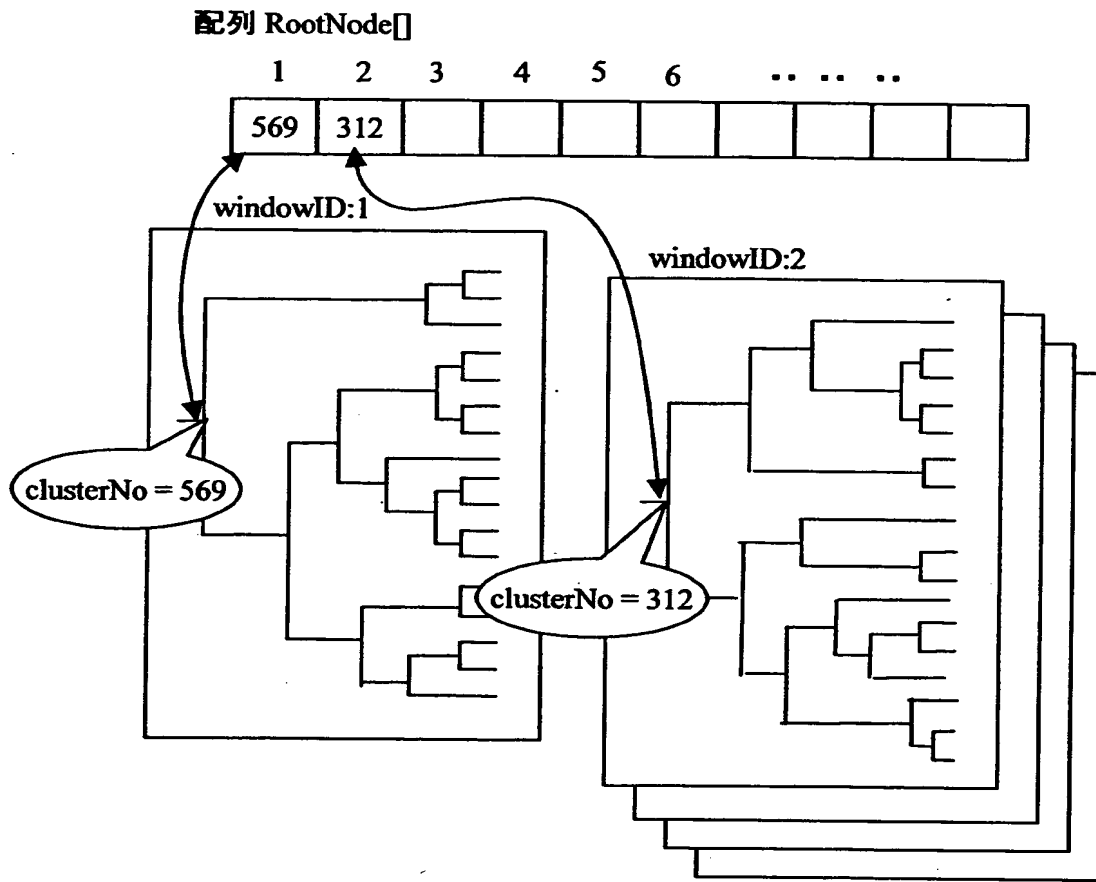
【図 1 3】



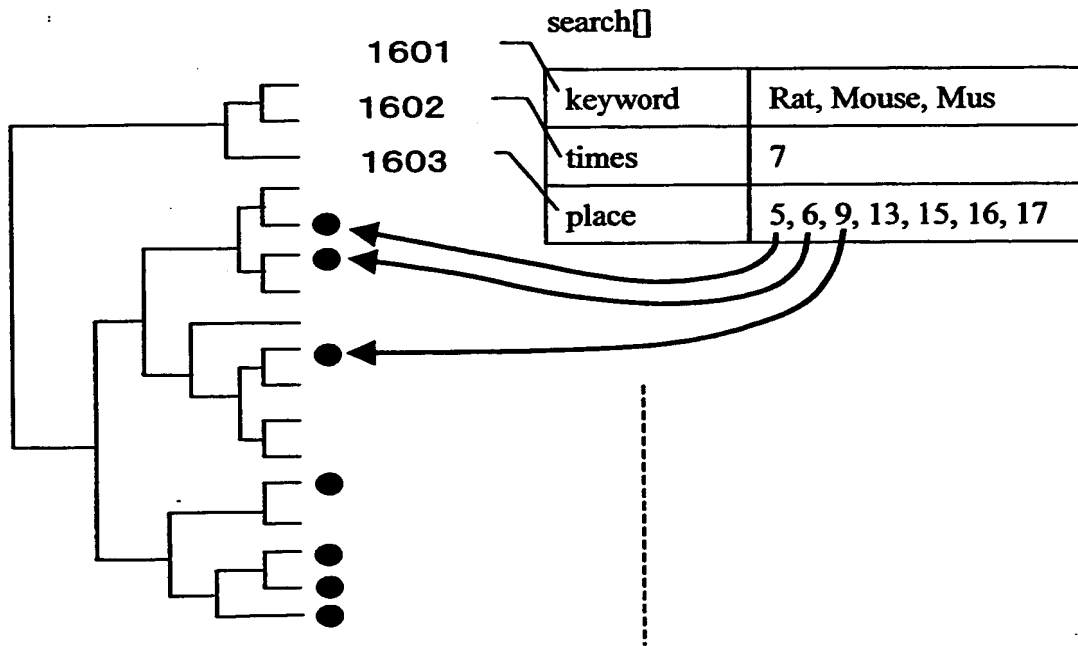
【図 1 4】



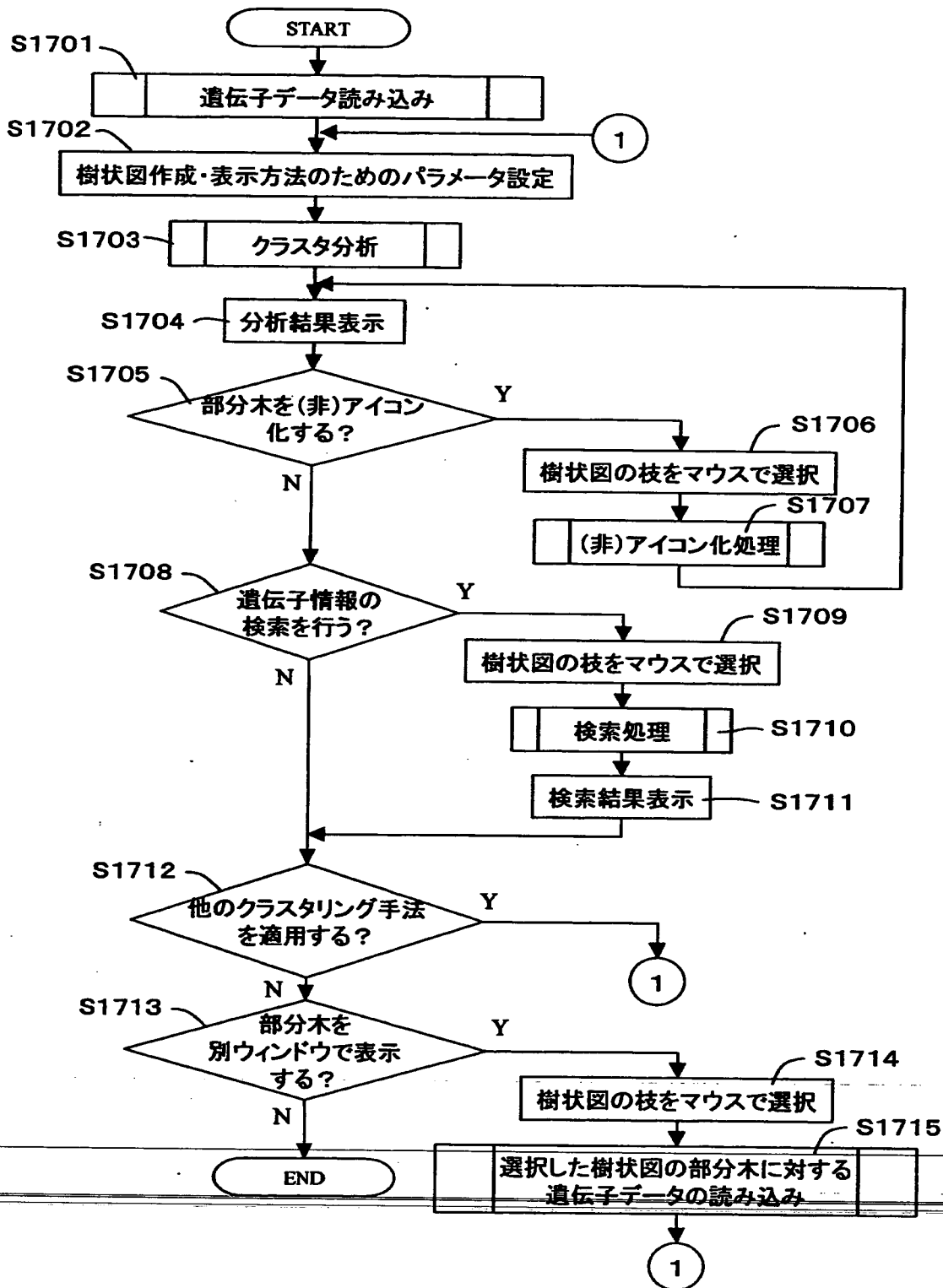
【図 1 5】



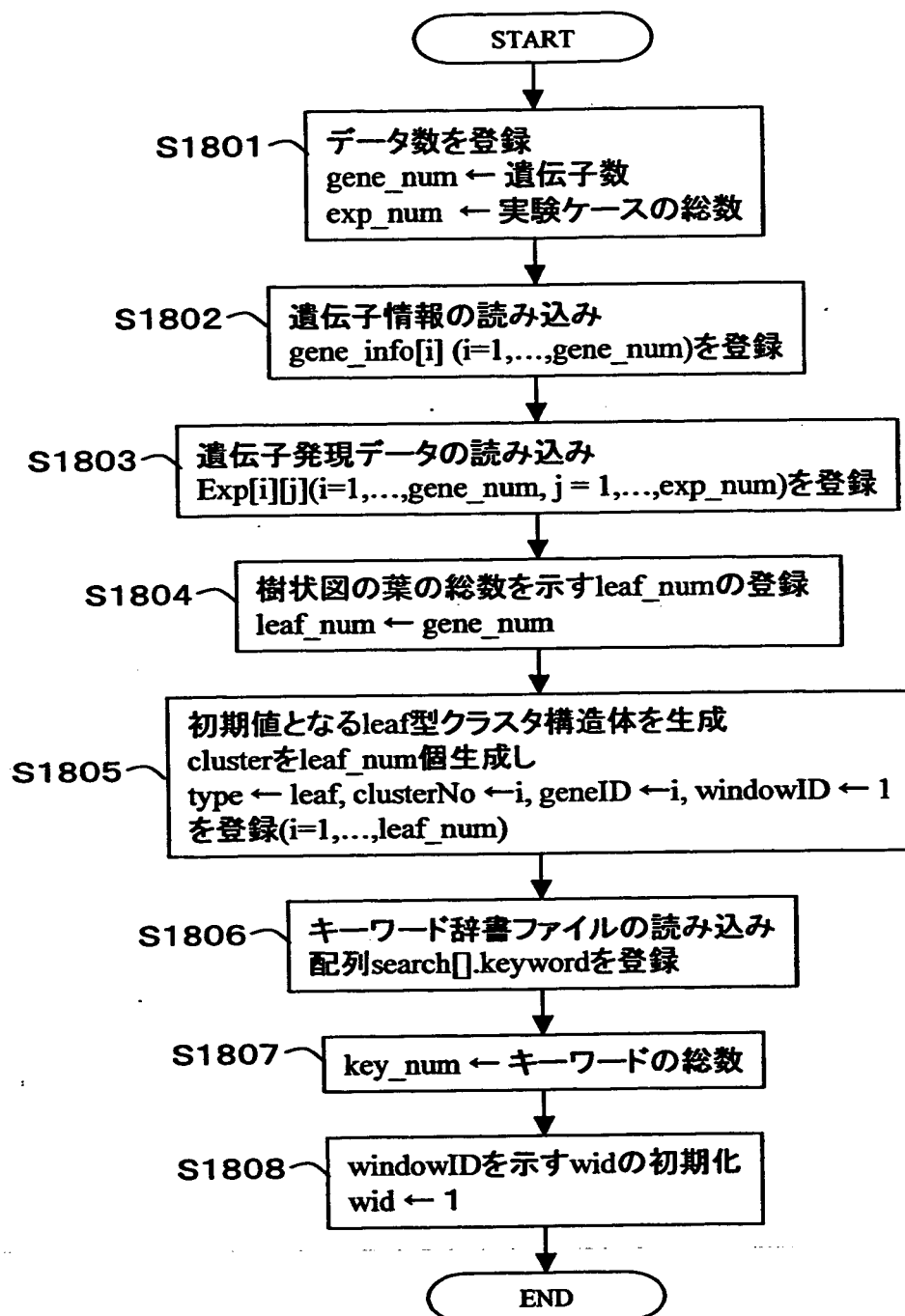
【図 1 6】



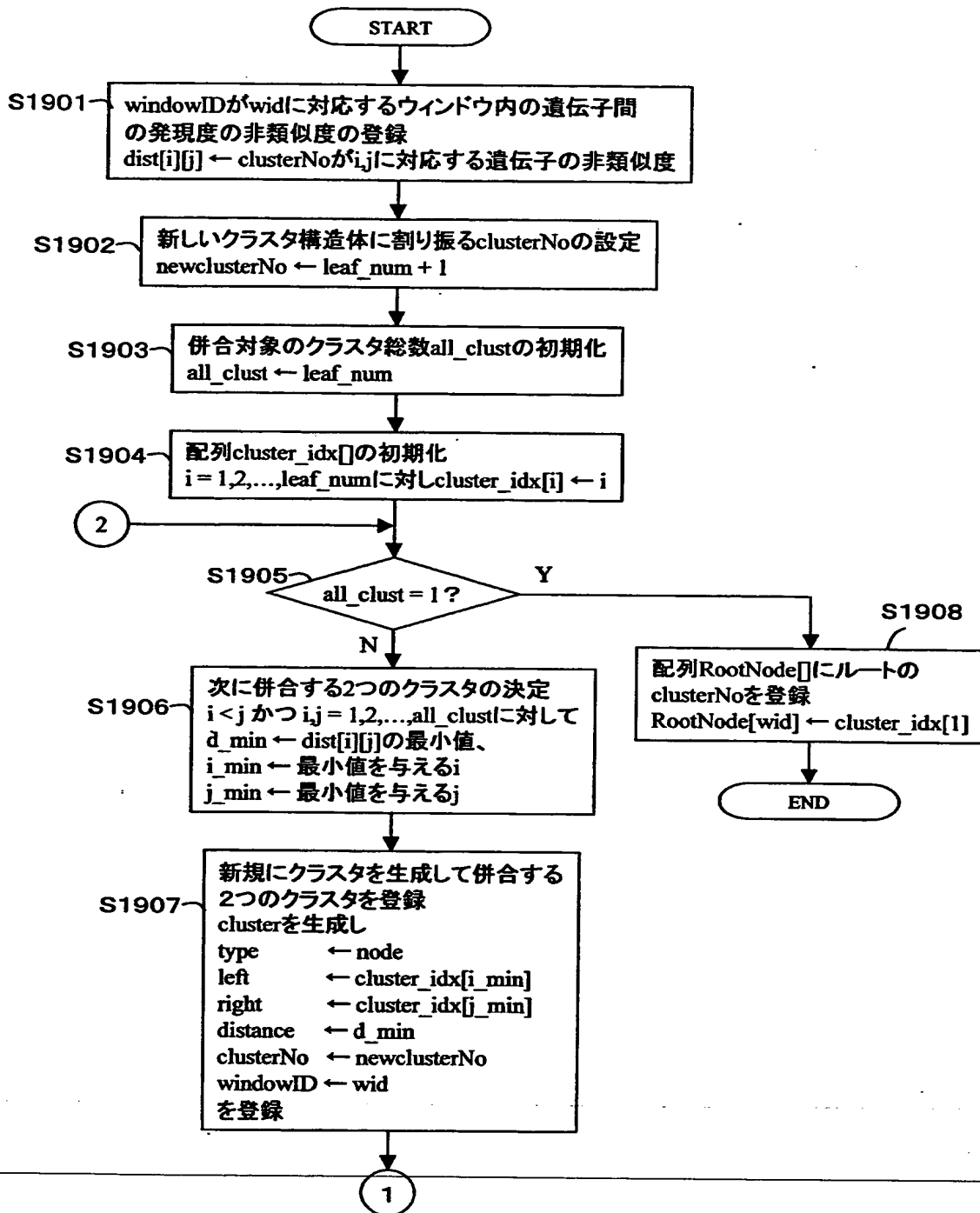
【図 17】



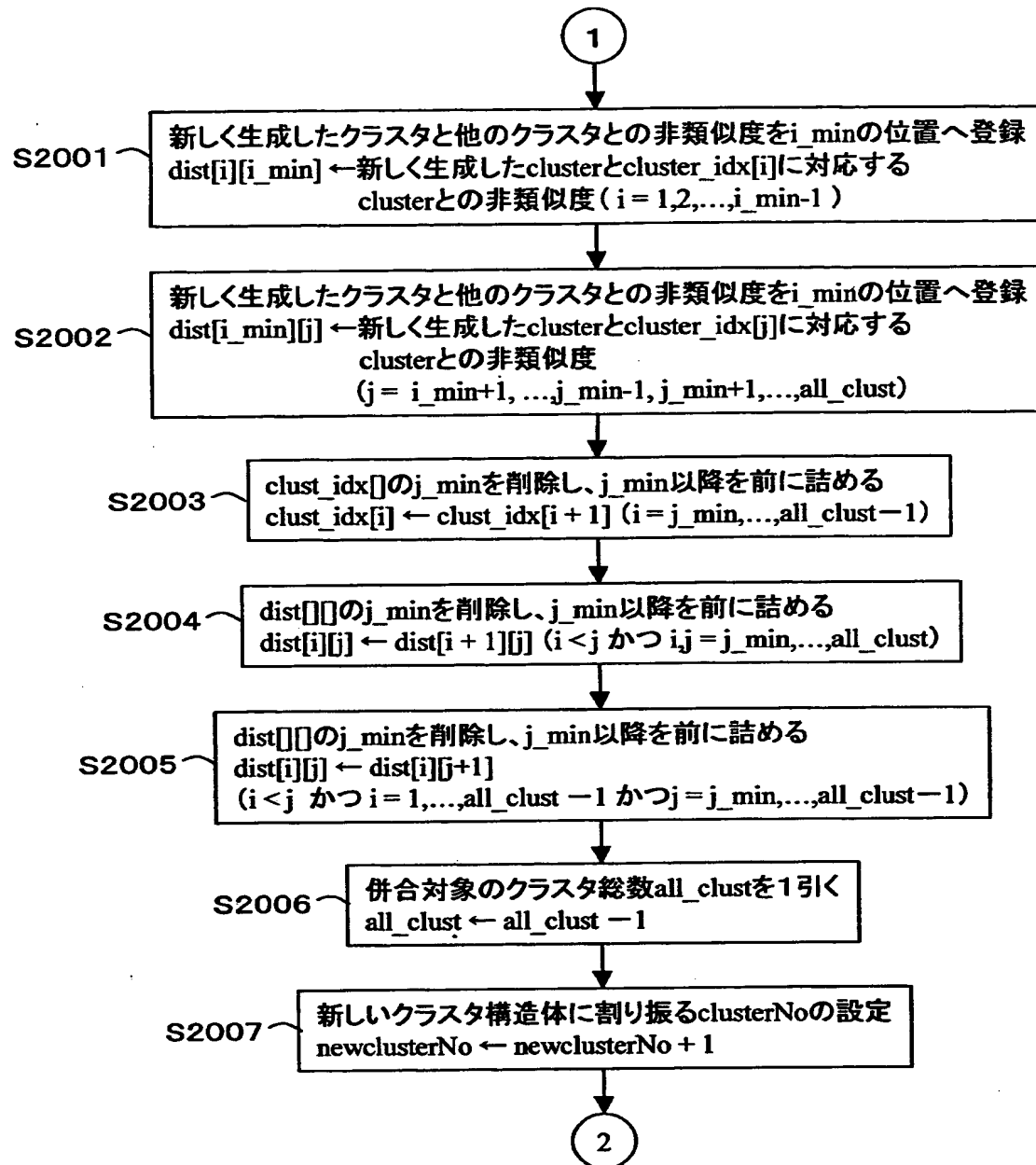
【図 18】



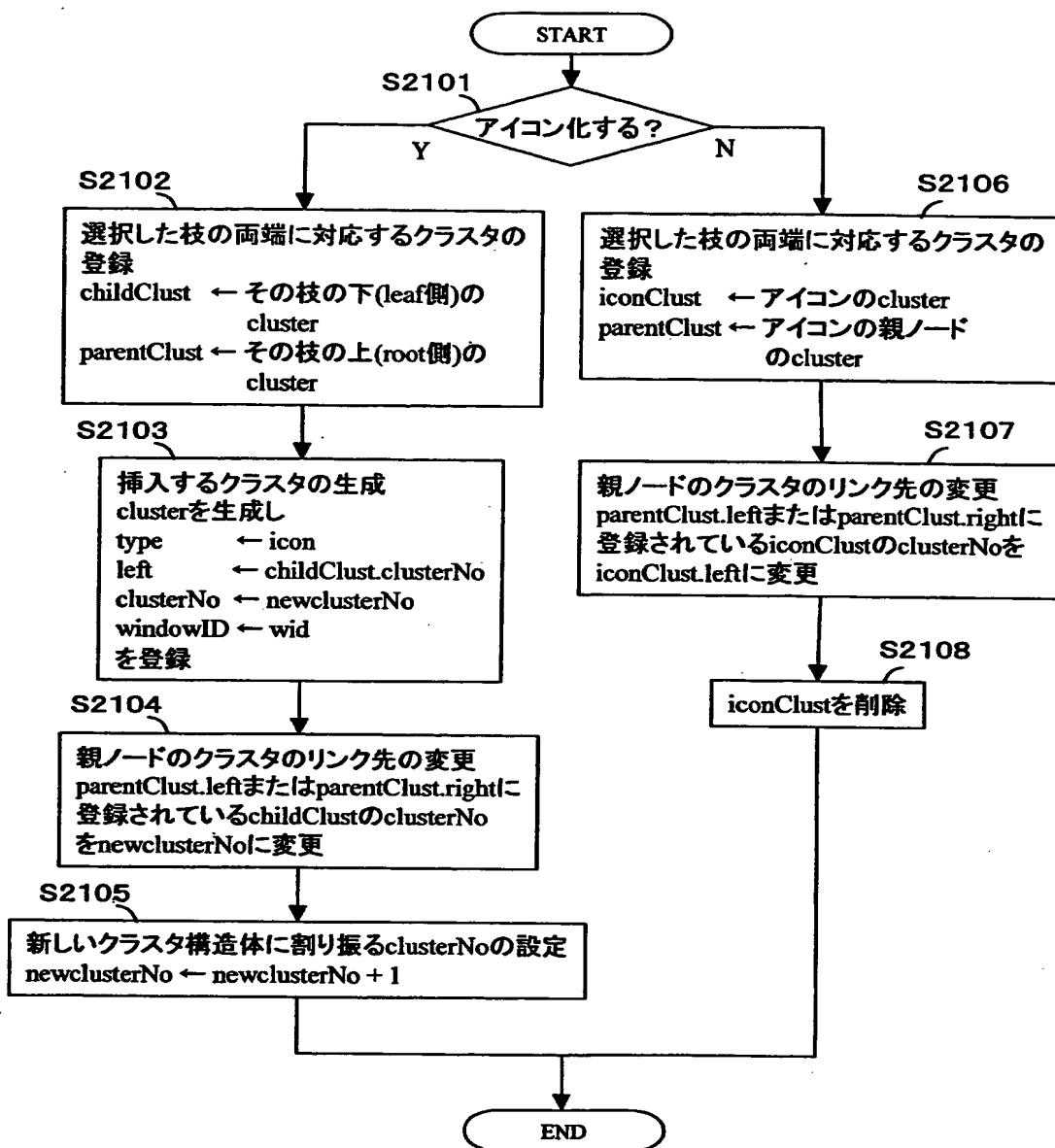
【図 1 9】



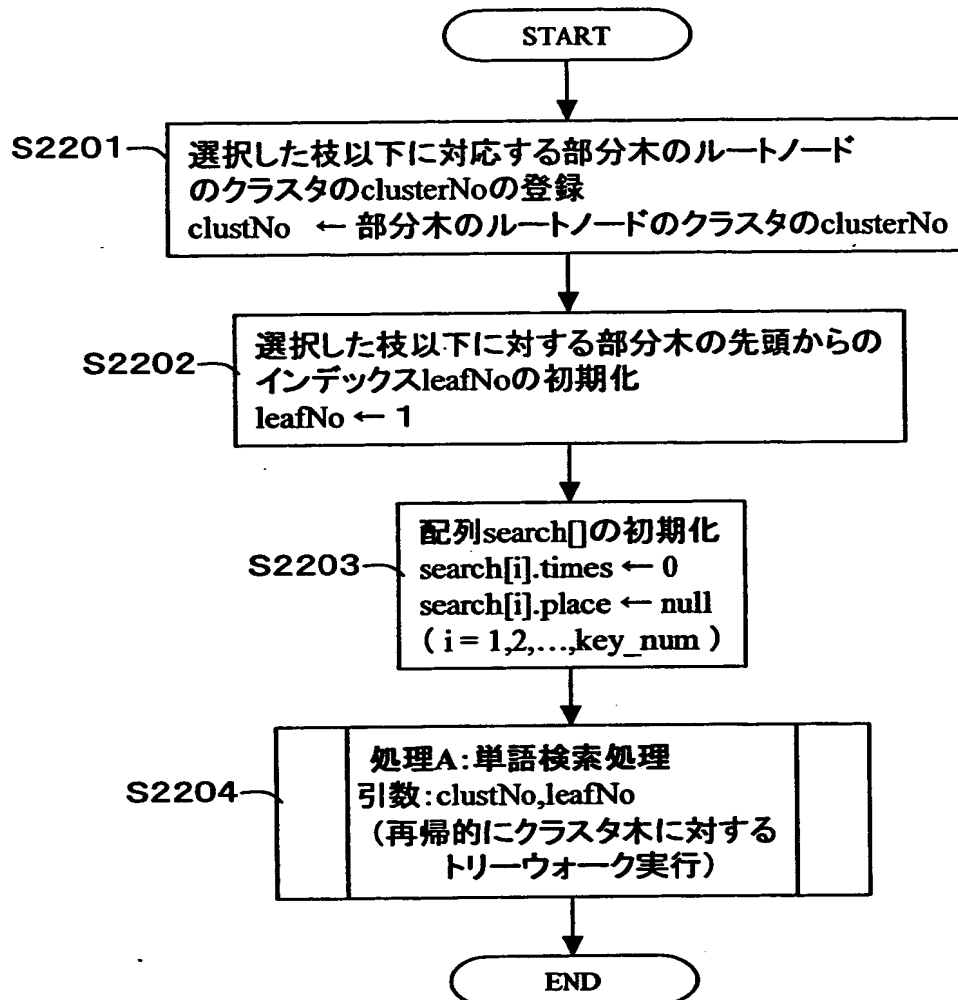
【図 2 0】



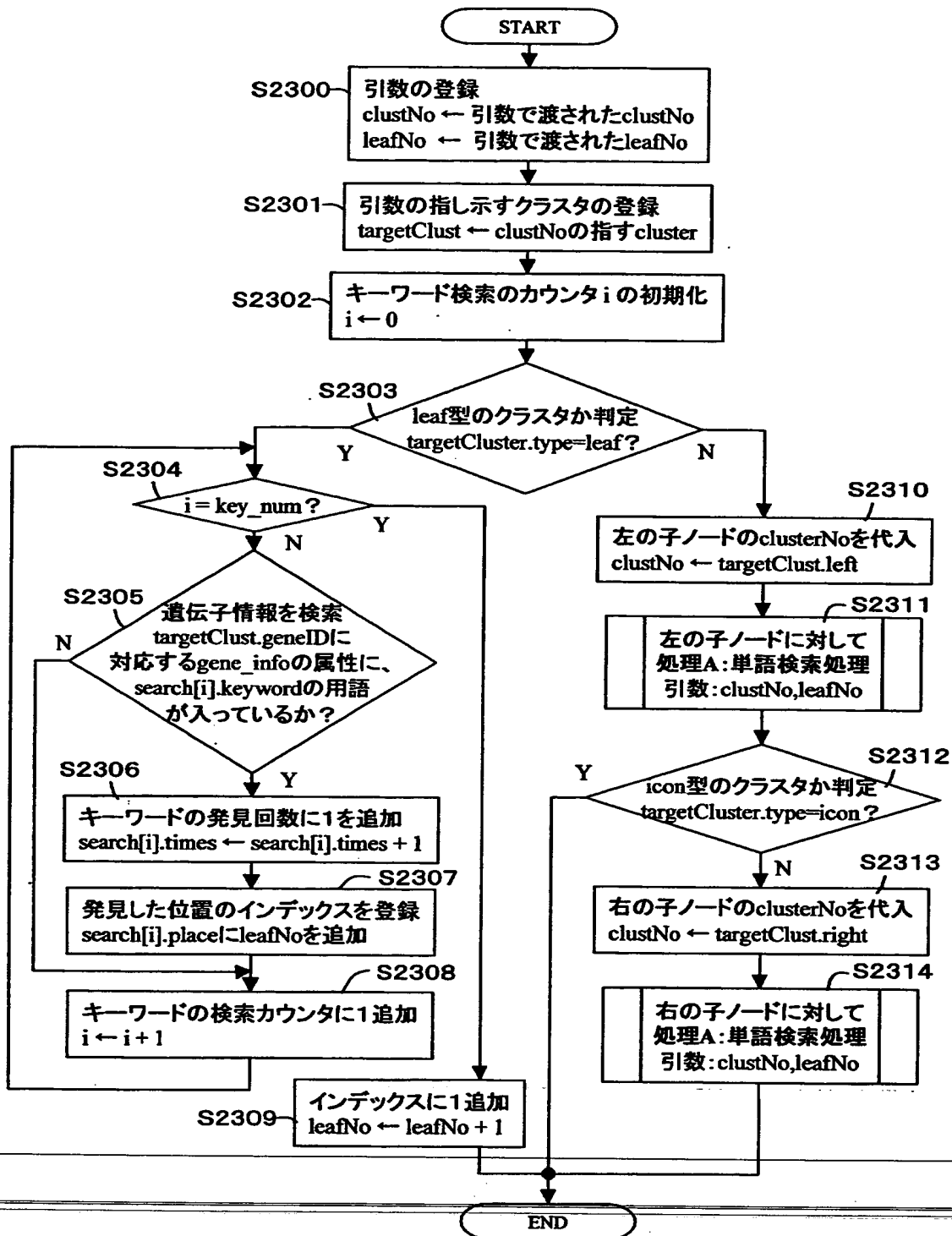
【図 21】



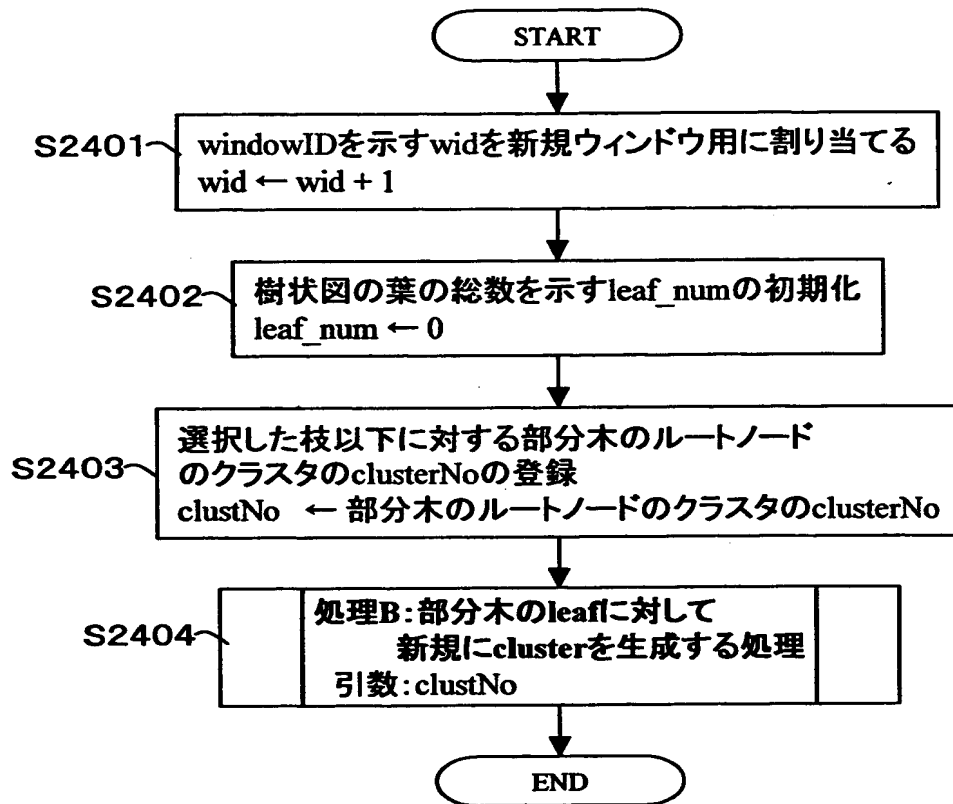
【図 22】



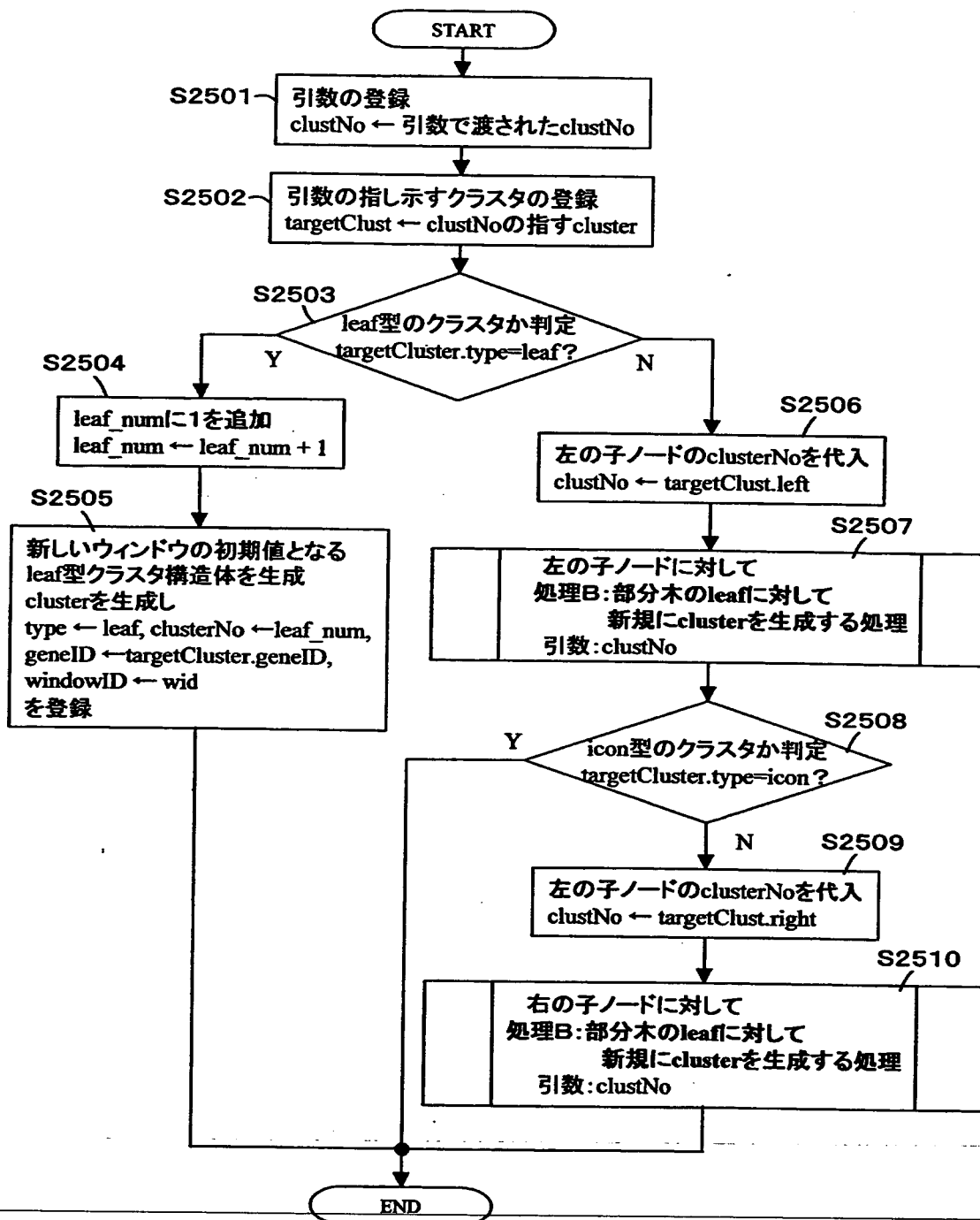
【図 2 3】



【図 2 4】



【図 2 5】



【書類名】 要約書

【要約】

【課題】 樹状図全体の枝の状態を大域的に把握すると共に個々の部分木の状態を詳細に知り、分類の絞り込みや、クラスタリング方法の選択の支援をする。

【解決手段】 樹状図の枝を選択し、選択した枝から葉の部分木に対して、別の表示ウィンドウで表示する機能、アイコン化する機能、アイコン化したものを元に戻す機能、部分木に含まれるキーワードを収集し表示する機能を備える。

【選択図】 図 5

出 願 人 履 歴 情 報

識別番号 [000233055]

1. 変更年月日 1990年 8月 7日

 [変更理由] 新規登録

 住 所 神奈川県横浜市中区尾上町6丁目81番地

 氏 名 日立ソフトウェアエンジニアリング株式会社

